



## Measuring and computing natural generators for homology groups

Chao Chen\*, Daniel Freedman

Rensselaer Polytechnic Institute, 110 8th street, Troy, NY 12180, USA

### ARTICLE INFO

#### Article history:

Received 4 June 2008

Received in revised form 6 April 2009

Accepted 16 June 2009

Available online 21 June 2009

Communicated by G. Rote

#### Keywords:

Computational topology

Computational geometry

Homology

Persistent homology

Homology basis

Stability

Finite field linear algebra

### ABSTRACT

We develop a method for measuring homology classes. This involves two problems. First, we define the size of a homology class, using ideas from relative homology. Second, we define an optimal basis of a homology group to be the basis whose elements' size have the minimal sum. We provide a greedy algorithm to compute the optimal basis and measure classes in it. The algorithm runs in  $O(\beta n^3 \log^2 n)$  time, where  $n$  is the size of the simplicial complex and  $\beta$  is the Betti number of the homology group. Finally, we prove the stability of our result. The algorithm can be adapted to measure any given class.

© 2009 Elsevier B.V. All rights reserved.

### 1. Introduction

The problem of computing the topological features of a space has recently drawn much attention from researchers in various fields, such as high-dimensional data analysis [1,2], graphics [3,4], networks [5] and computational biology [6,7]. Topological features are often preferable to purely geometric features, as they are more qualitative and global, and tend to be more robust. If the goal is to characterize a space, therefore, features which incorporate topology seem to be good candidates.

Once we are able to compute topological features, a natural problem is to rank the features according to their importance. The significance of this problem can be justified from two perspectives. First, unavoidable errors are introduced in data acquisition, in the form of traditional signal noise, and finite sampling of continuous spaces. These errors may lead to the presence of many small topological features that are not “real”, but are simply artifacts of noise or of sampling [8]. Second, many problems are naturally hierarchical. This hierarchy – which is a kind of multiscale or multi-resolution decomposition – implies that we want to capture the large scale features first. See Fig. 1 (Left, Center) for examples.

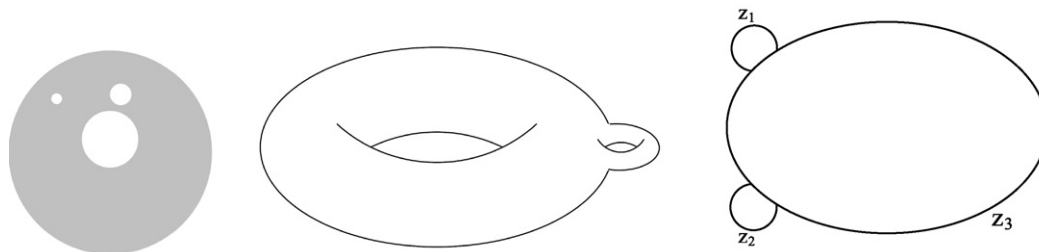
The topological features we use are homology groups over  $\mathbb{Z}_2$ , due to their ease of computation. (Thus, throughout this paper, all the additions are mod 2 additions.) We would then like to quantify or measure homology classes, as well as collections of classes. Specifically, there are two problems we would like to solve:

- (1) **Measuring the size of a homology class:** We need a way to quantify the size of a given homology class, and this size measure should agree with intuition. For example, in Fig. 1 (Left), the measure should be able to distinguish the one large class (of the 1-dimensional homology group) from the two smaller classes.

\* Corresponding author.

E-mail addresses: [chenc3@cs.rpi.edu](mailto:chenc3@cs.rpi.edu) (C. Chen), [freedd@cs.rpi.edu](mailto:freedd@cs.rpi.edu) (D. Freedman).

URLs: <http://www.cs.rpi.edu/~chenc3> (C. Chen), <http://www.cs.rpi.edu/~freedd> (D. Freedman).



**Fig. 1.** Left, Center: A disk with three holes and a 2-handled torus are really more like an annulus and a 1-handled torus, respectively, because the large features are more important. Right: A topological space formed from three circles. See accompanying discussion in the text.

(2) **Choosing a basis for a homology group:** We would like to choose a “good” set of homology classes to be the generators for the homology group (of a fixed dimension). Suppose that  $\beta$  is the dimension of this group, and that we are using  $\mathbb{Z}_2$  coefficients; then there are  $2^\beta - 1$  nontrivial homology classes in total. For a basis, we need to choose a subset of  $\beta$  of these classes, subject to the constraint that these  $\beta$  generate the group. The criterion of goodness for a basis is based on an overall size measure for the basis, which relies in turn on the size measure for its constituent classes. For instance, in Fig. 1 (Right), we must choose three from the seven nontrivial 1-dimensional homology classes:  $\{[z_1], [z_2], [z_3], [z_1] + [z_2], [z_1] + [z_3], [z_2] + [z_3], [z_1] + [z_2] + [z_3]\}$ . In this case, the intuitive choice is  $\{[z_1], [z_2], [z_3]\}$ , as this choice reflects the fact that there is really only one large cycle.

Furthermore, we make two additional requirements on the solution of aforementioned problems. First, the solution ought to be computable for topological spaces of arbitrary dimension. Second the solution should not require that the topological space be embedded, for example in a Euclidean space; and if the space is embedded, the solution should not make use of the embedding. These requirements are natural from the theoretical point of view, but may also be justified based on the following applications:

- In machine learning, it is often assumed that the data lives in a manifold whose dimension is much smaller than the dimension of the embedding space.
- In the study of shape, it is common to enrich the shape with other quantities, such as curvature, or color and other physical quantities. This leads to high-dimensional manifolds (e.g., 5–7 dimensions) embedded in high-dimensional ambient spaces [9].

Although there are existing low-dimensional techniques for approaching the problems we have laid out, to our knowledge, there are no definitions and algorithms satisfying the two requirements.

### 1.1. Related works

Persistent homology [10–15] is designed to track the lifetimes of homological features over the course of a filtration of a topological space. At first blush, it might seem that the powerful techniques of this theory are ideally suited to solving the problems we have set out. However, due to their somewhat different motivation, these techniques do not quite yield a solution. There are two reasons for this. First, the persistence of a feature depends not only on the space in which the feature lives, but also on the filtering function chosen. In the absence of a geometrically meaningful filter, it is not clear whether the persistence of a feature is a meaningful representation of its size. Second, and more importantly, the persistence only gives information for homology classes which ultimately die; for classes which are intrinsically part of the topological space, and which thus never die, the persistence is infinite. However, it is precisely these *essential* (or nonpersistent) classes that we care about. In more recent work, Cohen-Steiner et al. [16] have extended persistent homology in such a way that essential homology classes also have finite persistences. However, the persistences thus computed still depend on the filter function, and furthermore, do not always seem to agree with an intuitive notion of size.

Zomorodian and Carlsson [17] take a different approach to solving the localization problem. Their method starts with a topological space and a cover, a set of spaces whose union contains the original space. A blowup complex is built up which contains homology classes of all the spaces in the cover. The authors then use persistent homology to identify homology classes in the blowup complex which correspond to a same homology class in the given topological space. The persistent homology algorithm produces a complete set of generators for the relevant homology group, which forms a basis for the group. However, both the quality of the generators and the complexity of the algorithm depend strongly on the choice of cover; there is, as yet, no suggestion of a canonical cover.

Erickson and Whittlesey [18] showed how to compute the optimal basis for a 1-dimensional homology group in a 2-manifold. The authors also showed how the idea carries over to finding the optimal generators of the first fundamental group, though the proof is considerably harder in this case. A similar measure was used by Wood et al. [8] to remove topological noise of 2-dimensional surface. Both works are restricted to 2-dimensional topological space.

In this paper, we use the idea of growing geodesic balls. A similar idea has been used in [8,19]. However, this latter work depends on low-dimensional geometric reasoning, and hence is restricted to 1-dimensional homology classes in 2-manifold.

### 1.2. Our contributions

In this paper, we solve the aforementioned two problems. Our contributions include:

- Definitions of the size of homology classes and the optimal homology basis.
- A provably correct greedy algorithm to compute the optimal homology basis and measure its classes. This algorithm uses the persistent homology.
- An improvement of the straightforward algorithm using finite field linear algebra.
- A proof of the stability of our result with respect to small changes in certain quantities (to be explained in greater detail in Section 6).
- An algorithm to measure and localize a given class (Section 7).

## 2. Preliminaries

In this section, we briefly describe the background necessary for our work, including a discussion of homology groups, persistent homology and relative homology. Please refer to [20] for details of homology and relative homology, and [15] for persistent homology. For simplicity, we restrict our discussion to the combinatorial framework of simplicial homology over  $\mathbb{Z}_2$  field.

### 2.1. Homology groups

Within a given simplicial complex  $K$ , a  $d$ -chain is a formal sum of  $d$ -simplices in  $K$ ,  $c = \sum_{\sigma \in K} a_\sigma \sigma$ ,  $a_\sigma \in \mathbb{Z}_2$ . All the  $d$ -chains form the group of  $d$ -chains,  $C_d(K)$ . The boundary of a  $d$ -chain is the sum of the  $(d - 1)$ -faces of all the  $d$ -simplices in the chain. The boundary operator  $\partial_d : C_d(K) \rightarrow C_{d-1}(K)$  is a group homomorphism.

A  $d$ -cycle is a  $d$ -chain without boundary. The set of  $d$ -cycles forms a subgroup of the chain group, which is the kernel of the boundary operator,  $Z_d(K) = \ker(\partial_d)$ . A  $d$ -boundary is the boundary of a  $(d + 1)$ -chain. The set of  $d$ -boundaries forms a group, which is the image of the boundary operator,  $B_d(K) = \text{img}(\partial_{d+1})$ . It is not hard to see that a  $d$ -boundary is also a  $d$ -cycle. Therefore,  $B_d(K)$  is a subgroup of  $Z_d(K)$ . A  $d$ -cycle which is not a  $d$ -boundary,  $z \in Z_d(K) \setminus B_d(K)$ , is a nonbounding cycle. In our case, the coefficients belong to a field, namely  $\mathbb{Z}_2$ ; when this is the case, the groups of chains, boundaries and cycles are all vector spaces. Note that this is not true when the homology is over a ring which is not a field, such as  $\mathbb{Z}$ .

The  $d$ -dimensional homology group is defined as the quotient group  $H_d(K) = Z_d(K)/B_d(K)$ . An element in  $H_d(K)$  is a homology class, which is a coset of  $B_d(K)$ ,  $[z] = z + B_d(K)$  for some  $d$ -cycle  $z \in Z_d(K)$ . If  $z$  is a  $d$ -boundary,  $[z] = B_d(K)$  is the identity element of  $H_d(K)$ . Otherwise, when  $z$  is a nonbounding cycle,  $[z]$  is a nontrivial homology class and  $z$  is called a representative cycle of  $[z]$ . Cycles in the same homology class are homologous to each other, which means their difference is a boundary.

The dimension of the homology group, which is referred to as the Betti number,

$$\beta_d = \dim(H_d(K)) = \dim(Z_d(K)) - \dim(B_d(K)).$$

It can be computed with a reduction algorithm based on row and column operations of the boundary matrices [20]. Various reduction algorithms have been devised for different purposes [10,12,21].

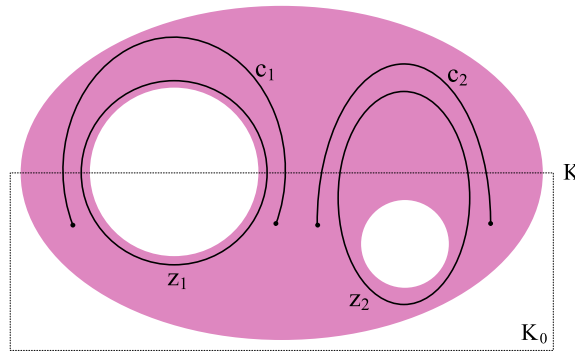
Note that since the field is  $\mathbb{Z}_2$ , the set of  $d$ -chains is in one-to-one correspondence with the set of subsets of  $d$ -simplices. A  $d$ -chain corresponds to a  $n_d$ -dimensional vector, whose nonzero entries correspond to the included  $d$ -simplices. Here  $n_d$  is the number of  $d$ -simplices in  $K$ . Computing the boundary of a  $d$ -chain corresponds to multiplying the chain vector with a boundary matrix  $[b_1, \dots, b_{n_d}]$ , whose column vectors are boundaries of  $d$ -simplices in  $K$ . By slightly abusing the notation, we call the boundary matrix  $\partial_d$ .

We call a subset of simplices of a given simplicial complex a subcomplex if this subset itself is a simplicial complex. The following notation will prove convenient. We say that a  $d$ -chain  $c \in C_d(K)$  is carried by a subcomplex  $K_0$  when all the  $d$ -simplices of  $c$  belong to  $K_0$ , formally,  $c \subseteq K_0$ . We denote  $\text{vert}(K)$  as the set of vertices of the simplicial complex  $K$ ,  $\text{vert}(c)$  as that of the chain  $c$ .

Replacing simplexes by their continuous images in a given topological space gives singular homology. The simplicial homology of a simplicial complex is naturally isomorphic to the singular homology of its geometric realization. This implies, in particular, that the simplicial homology of a space does not depend on the particular simplicial complex chosen for the space. In figures of this paper, we often ignore the simplicial complex and only show the continuous images of chains.

### 2.2. Persistent homology

Given a topological space  $\mathbb{X}$  and a filter function  $f : \mathbb{X} \rightarrow \mathbb{R}$ , persistent homology studies the homology classes of the sublevel sets,  $\mathbb{X}^t = f^{-1}(-\infty, t]$ . A nontrivial homology class in  $\mathbb{X}^{t_1}$  may become trivial in  $\mathbb{X}^{t_2}$ ,  $t_1 < t_2$  (formally, when



**Fig. 2.** A disk with two holes, whose triangulation is  $K$ . Simplices of  $K$  lying completely in the dotted rectangle form a subcomplex  $K_0$ . The 1-dimensional relative homology group  $H_1(K, K_0)$  has dimension 1, although  $H_1(K)$  has dimension 2. The nontrivial class  $[z_2]$  is carried by  $K_0$ .

induced by the inclusion homomorphism). Persistent homology tries to capture this phenomenon by measuring the times at which a homology class is born and dies. The persistence, or life time of the class is the difference between its death and birth times. Those with longer lives tell us something about the global structure of the space  $\mathbb{X}$ , as described by the filter function. Note that the *essential*, that is, nontrivial homology classes of the given topological space  $\mathbb{X}$  will never die.

Edelsbrunner et al. [10] devised an  $O(n^3)$  algorithm to compute the persistent homology. Its input are a simplicial complex  $K$  and a filter function  $f$ , which assigns each simplex in  $K$  a real value. Simplices of  $K$  are sorted in ascending order according to their filter function values. This order is actually the order in which simplices enter the sublevel set  $f^{-1}(-\infty, t]$  while  $t$  increases. For simplicity, in this paper we call this ordering the *simplex-ordering* of  $K$  with regard to  $f$ . Note that within the simplex-ordering, a simplex must appear after all of its faces. With this restriction, any sublevel set is a subcomplex. All the sublevel sets taken together form a *filtration*, namely, a nested sequence that begins with the empty complex and ends with the complete complex,  $\emptyset = K_0 \subset K_1 \subset \dots \subset K_m = K$ . Given this input, the output of the algorithm is the birth and death times of homology classes.

More specifically, the persistence algorithm is based on column reductions of boundary matrices. The latest version of this algorithm [7,15] unifies boundary matrices of different dimensions into one overall *incidence matrix*  $D$ . Rows and columns of  $D$  correspond to simplices of  $K$ , indexed in the simplex-ordering. An entry of  $D$  is 1 if and only if its corresponding entry is 1 in the corresponding boundary matrix. The algorithm performs column reductions on  $D$  from left to right. Each new column is reduced by addition with the already reduced columns, until its lowest nonzero entry is as high as possible. The reduced matrix  $R = DV$  provides  $\text{rank}(D)$  pairings of simplices, in which each simplex appears at most once. The filter function values of each pairing are the birth and death times of a persistent homology class. Unpaired simplices are paired with  $+\infty$  and correspond to essential classes. Simplices paired with  $+\infty$  or paired on the left are *positive*, and the rest are *negative*.

The reduction is completely recorded in the matrix  $V$ . Columns of  $V$  corresponding to positive simplices form bases of cycle groups. Columns corresponding to positive simplices paired with  $+\infty$  are cycles representing essential classes.

### 2.3. Relative homology

Given a simplicial complex  $K$  and a subcomplex  $K_0 \subseteq K$ , we may wish to study the structure of  $K$  by ignoring all the chains in  $K_0$ . We consider two  $d$ -chains,  $c_1$  and  $c_2$  to be the same if their difference is carried by  $K_0$ . The objects we are interested in are then defined as these equivalence classes, which form a quotient group,  $C_d(K, K_0) = C_d(K)/C_d(K_0)$ . We call it the *group of relative chains*, whose elements (cosets), are called *relative chains*.

The boundary operator  $\partial_d : C_d(K) \rightarrow C_{d-1}(K)$  induces a *relative boundary operator*,  $\partial_d^{K_0} : C_d(K, K_0) \rightarrow C_{d-1}(K, K_0)$ . Analogous to the way we define  $Z_d(K)$ ,  $B_d(K)$  and  $H_d(K)$  in  $C_d(K)$ , we define the *group of relative cycles*, the *group of relative boundaries* and the *relative homology group* in  $C_d(K, K_0)$ , denoted as  $Z_d(K, K_0)$ ,  $B_d(K, K_0)$  and  $H_d(K, K_0)$ , respectively. An element in  $Z_d(K, K_0) \setminus B_d(K, K_0)$  is a *nonbounding relative cycle*.

The following notation will prove convenient. We define a homomorphism  $\phi_{K_0} : C_d(K) \rightarrow C_d(K, K_0)$  mapping  $d$ -chains to their corresponding relative chains,  $\phi_{K_0}(c) = c + C_d(K_0)$ . This homomorphism induces another homomorphism,  $\phi_{K_0}^* : H_d(K) \rightarrow H_d(K, K_0)$ , mapping homology classes of  $K$  to their corresponding relative homology classes,  $\phi_{K_0}^*(h) = \phi_{K_0}(z) + B_d(K, K_0)$  for any  $z \in h$ .

Given a  $d$ -chain  $c \in C_d$ , its corresponding relative chain  $\phi_{K_0}(c)$  is a relative cycle if and only if  $\partial_d(c)$  is carried by  $K_0$ . Furthermore, it is a relative boundary if and only if there is a  $(d + 1)$ -chain  $c' \in C_{d+1}(K)$  such that  $c - \partial_{d+1}(c')$  is carried by  $K_0$ .

These ideas are illustrated in Fig. 2. Although  $z_1$  and  $z_2$  are both nonbounding cycles in  $K$ ,  $\phi_{K_0}(z_1)$  is a nonbounding relative cycle whereas  $\phi_{K_0}(z_2)$  is only a relative boundary. Although chains  $c_1$  and  $c_2$  are not cycles in  $K$ ,  $\phi_{K_0}(c_1)$  and  $\phi_{K_0}(c_2)$  are relative cycles homologous to  $\phi_{K_0}(z_1)$  and  $\phi_{K_0}(z_2)$ , respectively.

Note that  $[z_1]$  and  $[z_2]$  are both nontrivial homology classes in  $K$ . But their corresponding classes in the relative homology group may be trivial. We say a subcomplex  $K_0$  carries a class  $h$  if  $h$  has a trivial image in the relative homology group  $H_d(K, K_0)$ , formally,  $\phi_{K_0}^*(h) = 0 + B_d(K, K_0)$ . Intuitively, this means that  $h$  disappears if we delete  $K_0$  from  $K$ , by contracting it into a point and modding it out. The following lemma gives us more intuition behind this definition.

**Lemma 1.**  $K_0$  carries  $h$  if and only if it carries a cycle of  $h$ .

**Proof.** For any cycle  $z \in h$ , the relative chain  $\phi_{K_0}(z)$  is a relative boundary if and only if there is a  $(d + 1)$ -chain  $c' \in C_{d+1}(K)$  such that  $z - \partial_{d+1}(c') \in h$  is carried by  $K_0$ .  $\square$

For example, in Fig. 2,  $\phi_{K_0}^*([z_1])$  is a nontrivial relative homology class, whereas  $\phi_{K_0}^*([z_2])$  is trivial. We say that the class  $[z_2]$  is carried by  $K_0$ . This concept plays an important role in our definition of the size measure. Further details will be given in Section 3.2.

### 2.4. Rank computations of sparse matrices over finite fields

Wiedemann [22] presented a randomized algorithm to capture the rank of a sparse matrix over finite field. The expected time of the algorithm is  $O(n(\omega + n \log n) \log n)$ , where  $n$  is the maximal dimension of the matrix and  $\omega$  is the total number of nonzero entries.

## 3. Defining the problem

In this section, we provide a technique for ranking homology classes according to their importance. Specifically, we solve the two problems mentioned in Section 1 by formally defining (1) a meaningful size measure for homology classes that is computable in arbitrary dimension; and (2) an optimal homology basis which distinguishes large classes from small ones effectively.

### 3.1. The discrete geodesic distance

In order to measure the size of homology classes, we need a notion of distance. As we will deal with a simplicial complex  $K$ , it is most natural to introduce a discrete metric, and corresponding distance functions. We define the *discrete geodesic distance* from a vertex  $p \in \text{vert}(K)$ ,  $f_p : \text{vert}(K) \rightarrow \mathbb{R}$ , as follows. Suppose each edge in  $K$  has a nonnegative weight, for any vertex  $q \in \text{vert}(K)$ ,  $f_p(q) = \text{dist}(p, q)$  is the length of the shortest path connecting  $p$  and  $q$ , in the 1-skeleton of  $K$ . We may then extend this distance function from vertices to higher-dimensional simplices naturally. For any simplex  $\sigma \in K$ ,  $f_p(\sigma)$  is the maximal function value of the vertices of  $\sigma$ ,  $f_p(\sigma) = \max_{q \in \text{vert}(\sigma)} f_p(q)$ . This extension has the same effect as linearly interpolating the function on the interiors of the simplices (the sublevel sets of the two extensions are homotopy equivalent). Finally, we define a geodesic ball  $B_p^r$ ,  $p \in \text{vert}(K)$ ,  $r \geq 0$ , as the subset of  $K$ ,  $B_p^r = \{\sigma \in K \mid f_p(\sigma) \leq r\}$ . It is straightforward to show that these subsets are in fact subcomplexes.

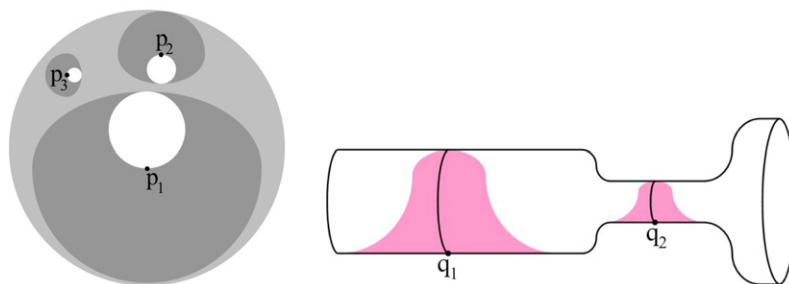
### 3.2. Measuring the size of a homology class

Using relative homology, we define the size of a homology class as follows. Given a simplicial complex  $K$ , assume we are given a collection of subcomplexes  $\mathcal{L} = \{L \subseteq K\}$ . Furthermore, each of these subcomplexes is endowed with a size. In this case, we define the size of a homology class  $h$  as the size of the smallest  $L$  carrying  $h$  (assuming one such  $L$  exists, which can be guaranteed if  $\mathcal{L}$  is properly chosen).

**Definition 2.** The size of a class  $h$ ,  $S(h)$ , is the size of the smallest measurable subcomplex carrying  $h$ , formally,

$$S(h) = \min_{L \in \mathcal{L}} \text{size}(L) \quad \text{s.t.} \quad \phi_L^*(h) = B_d(K, L).$$

In this paper, we take  $\mathcal{L}$  to be the set of discrete geodesic balls,  $\mathcal{L} = \{B_p^r \mid p \in \text{vert}(K), r \geq 0\}$ . The size of a geodesic ball is naturally its radius  $r$ . The smallest geodesic ball carrying  $h$  is denoted as  $B_{\min}(h)$  for convenience, whose radius is  $S(h)$ . In Fig. 3 (Left), the three geodesic balls centered at  $p_1$ ,  $p_2$  and  $p_3$  are the smallest geodesic balls carrying nontrivial homology classes corresponding to the three holes. Their radii are the size of the three classes. In Fig. 3 (Right), the smallest geodesic ball carrying a nontrivial homology class is the shaded one centered at  $q_2$ , not the one centered at  $q_1$ . Note that these geodesic ball may not look like Euclidean balls in the embedding space.



**Fig. 3.** Left: On a disk with three holes, the three shaded regions are the three smallest geodesic balls measuring the three corresponding classes. Right: On a tube, the smallest geodesic ball is centered at  $q_2$ , not  $q_1$ .

### 3.3. The optimal homology basis

For the  $d$ -dimensional  $\mathbb{Z}_2$  homology group whose dimension (Betti number) is  $\beta_d$ , there are  $2^{\beta_d} - 1$  nontrivial homology classes. However, we only need  $\beta_d$  of them to form a basis. The basis should be chosen wisely so that we can easily distinguish important homology classes from noise. See Fig. 1 (Right) for an example. There are  $2^3 - 1 = 7$  nontrivial homology classes; we need three of them to form a basis. We would prefer to choose  $\{[z_1], [z_2], [z_3]\}$  as a basis, rather than  $\{[z_1] + [z_2] + [z_3], [z_2] + [z_3], [z_3]\}$ . The former indicates that there is one big cycle in the topological space, whereas the latter gives the impression of three large classes.

In keeping with this intuition, the *optimal homology basis* is defined as follows.

**Definition 3.** The optimal homology basis is the basis for the homology group whose elements' size have the minimal sum, formally,

$$\mathcal{H}_d = \operatorname{argmin}_{\{h_1, \dots, h_{\beta_d}\}} \sum_{i=1}^{\beta_d} S(h_i) \quad \text{s.t.} \quad \dim(\{h_1, \dots, h_{\beta_d}\}) = \beta_d.$$

This definition guarantees that large homology classes appear as few times as possible in the optimal homology basis. In Fig. 1 (Right), the optimal basis will be  $\{[z_1], [z_2], [z_3]\}$ , which has only one large class.

This definition uses  $L_1$ -norm on the vector of sizes. Since all class sizes are nonnegative, and further, since the problem has a matroid structure (to be demonstrated in the next section), it will turn out that we can use any  $L_p$ -norm in the definition and still get the same optimal homology basis. An exception, however, is the  $L_\infty$ -norm. In this case, there may be many different optimal bases. The optimal basis defined using  $L_1$ -norm is one of them.

For each class in the basis, we need a cycle representing it. According to Lemma 1,  $B_{\min}(h)$ , the smallest geodesic ball carrying  $h$ , carries at least one cycle of  $h$ . We localize each class in the optimal basis by its *localized cycles*, which are cycles of  $h$  carried by  $B_{\min}(h)$ . This is a fair choice because it is consistent with the size measure of  $h$  and it is computable in polynomial time.

Please note that the localized cycle may not be the simplest one. The cycle may wiggle a lot inside the geodesic ball,  $B_{\min}(h)$ . The authors addressed this issue in much greater detail in [23,24]. In these papers, different size definitions are provided; for example, the localized cycle may be defined as the representative cycle with the minimal number of simplices. These new definitions give representative cycles which are simple (concise) in both geometry and algebra. Unfortunately (and perhaps not surprisingly), it turns they are NP-hard to compute and even to approximate.

## 4. The algorithm

In this section, we introduce an algorithm to compute the optimal homology basis as defined in Definition 3. For each class in the basis, we measure its size, and represent it with one of its localized cycles. We first introduce an algorithm to compute the smallest homology class, namely,  $\text{Measure-Smallest}(K)$ . Based on this procedure, we provide the algorithm  $\text{Measure-All}(K)$ , which computes the optimal homology basis. The algorithm takes  $O(\beta_d n^4)$  time, where  $\beta_d$  is the Betti number for  $d$ -dimensional homology group and  $n$  is the cardinality of the input simplicial complex  $K$ . Please note that in the rest of the paper, we assume  $d$ , the dimension of the relevant homology group, is given.

### 4.1. Computing the smallest homology class

The procedure  $\text{Measure-Smallest}(K)$  measures and localizes the smallest nontrivial homology class, namely, the one with the smallest size,

$$h_{\min} = \operatorname{argmin}_{h \in H_d(K): h \neq B_d(K)} S(h).$$

The output of this procedure will be a pair  $(S_{\min}, z_{\min})$ , namely, the size and a localized cycle of  $h_{\min}$ . According to the definitions, this pair is determined by the smallest geodesic ball carrying  $h_{\min}$ , namely,  $B_{\min}(h_{\min})$ .

It is straightforward to see that the ball  $B_{\min}(h_{\min})$  is also the smallest geodesic ball carrying any nontrivial homology class of  $K$ . It can be computed by computing  $B_p^{r(p)}$  for all vertices  $p$ , where  $B_p^{r(p)}$  is the smallest geodesic ball centered at  $p$  which carries any nontrivial homology class. When all the  $B_p^{r(p)}$ 's are computed, we compare their radii,  $r(p)$ 's, and pick the smallest ball as  $B_{\min}(h_{\min})$ .

For each vertex  $p$ , we compute  $B_p^{r(p)}$  by applying the persistent homology algorithm to  $K$  with the discrete geodesic distance from  $p$ ,  $f_p$ , as the filter function. Note that a geodesic ball  $B_p^r$  is the sublevel set  $f_p^{-1}(-\infty, r] \subseteq K$ . Nontrivial homology classes of  $K$  are essential homology classes in the persistent homology algorithm. (In the rest of this paper, we may use “essential homology classes” and “nontrivial homology classes of  $K$ ” interchangeably.) Therefore, the birth time of the first essential homology class, namely, the filter function value of the very first  $d$ -simplex that is positive and paired with  $+\infty$ , is  $r(p)$ , and the subcomplex  $f_p^{-1}(-\infty, r(p)]$  is  $B_p^{r(p)}$ .

Once we determine  $B_{\min}(h_{\min})$ , whose center is denoted as  $p_{\min}$ , the size  $S_{\min}$  is the radius  $r(p_{\min})$ . A localized cycle can be decided by the persistent homology algorithm with the filter function  $f_{p_{\min}}$ . Recall that the matrix  $V$  in the persistence reduction  $R = DV$  provides cycles representing essential classes. The localized cycle  $z_{\min}$  is the column of  $V$  corresponding to the very first  $d$ -simplex that is positive and paired with  $+\infty$ . The reason is it represents the youngest essential class, which is  $h_{\min}$ . Plus it is carried by  $B_{\min}(h_{\min})$ .

#### 4.2. Computing the optimal homology basis

In this section, we present the algorithm for computing the optimal homology basis defined in Definition 3, namely,  $\mathcal{H}_d$ . We first show that the optimal homology basis can be computed in a greedy manner. Second, we introduce an efficient greedy algorithm.

##### 4.2.1. Computing $\mathcal{H}_d$ in a greedy manner

As has been noted, over the  $\mathbb{Z}_2$  field, the homology group is a vector space. It, together with the family of its linearly independent subsets, form a vector matroid. Using the size of homology classes as a weight function, we have a weighted matroid. Matroid theory [25,26] suggests a greedy method to compute the optimal homology basis as follows.

For convenience, let  $H$  be the set of nontrivial  $d$ -dimensional homology classes (i.e. the homology group minus the trivial class). Denote  $\operatorname{seq}(H) = (h_1, h_2, \dots, h_{(2^{\beta_d}-1)})$  as the sequence of all the classes of  $H$  sorted in the monotonically increasing order according to size, formally,  $S(h_i) \leq S(h_{i+1})$  for all  $i$ . Repeatedly compute the smallest class in  $\operatorname{seq}(H)$  and pick each one which is linearly independent of those we have already picked, until  $\beta_d$  are picked. The collected  $\beta_d$  classes  $\{h_{i_1}, h_{i_2}, \dots, h_{i_{\beta_d}}\}$  form the optimal homology basis  $\mathcal{H}_d$ . (Note that the  $h$ 's are ordered by size, i.e.  $S(h_{i_k}) \leq S(h_{i_{k+1}})$ .)

However, this naive method may be exponential in  $\beta_d$ . For example, we may have to compute all the linear combinations of  $\{h_{i_1}, h_{i_2}, \dots, h_{i_{(\beta_d-1)}}\}$  before we find  $h_{i_{\beta_d}}$ . Next, we present our greedy algorithm which is polynomial.

##### 4.2.2. Computing $\mathcal{H}_d$ with a sealing technique

In this section, we introduce a polynomial greedy algorithm for computing  $\mathcal{H}_d$ . Instead of computing the smallest classes in  $\operatorname{seq}(H)$  one by one, our algorithm uses a sealing technique and takes time polynomial in  $\beta_d$  and  $n$ . Intuitively, when the smallest  $l$  classes in  $\mathcal{H}_d$  are picked, we make them trivial by adding new cells to the given complex. In the augmented complex, any linear combination of these picked classes becomes trivial, and the smallest nontrivial class is the  $(l + 1)$ -th one in  $\mathcal{H}_d$ .

The algorithm starts by measuring and localizing the smallest homology class of the given simplicial complex  $K$  (using the procedure Measure-Smallest( $K$ ) introduced in Section 4.1), which is also the first class we choose for  $\mathcal{H}_d$ . We make this class trivial by sealing one of its cycles – i.e. the localized cycle we computed – with a new cell. Next, we measure and localize the smallest homology class of the augmented complex  $K'$ . This class is the second smallest homology class in  $\mathcal{H}_d$ . We make this class trivial again and proceed for the third smallest class in  $\mathcal{H}_d$ . This process is repeated for  $\beta_d$  rounds, yielding  $\mathcal{H}_d$ .

We make a homology class trivial by sealing the class's localized cycle, which we have computed. To seal this cycle  $z$ , we add a new  $(d + 1)$ -cell whose boundary is exactly this cycle. In Fig. 4, a 1-cycle with four edges,  $z_1$ , is sealed up with one new 2-cell. Please note that the new cell is not a simplex and the augmented complex  $K'$  is a cell complex, not a simplicial complex.

It is essential to make sure the new cell does not influence our measurement. We assign the new cell  $+\infty$  filter function values, formally,  $f_p(\sigma) = +\infty$  for all  $p \in \operatorname{vert}(K)$  and  $\sigma \in K' \setminus K$ .

The algorithm is illustrated in Fig. 4. Assuming unit edge lengths, the 4-edge cycle,  $z_1$ , and the 8-edge cycle,  $z_2$ , are the localized cycles of the smallest and the second smallest homology classes ( $S([z_1]) = 2$ ,  $S([z_2]) = 4$ ). The nonbounding cycle  $z_3 = z_1 + z_2$  corresponds to the largest nontrivial homology class  $[z_3] = [z_1] + [z_2]$  ( $S([z_3]) = 5$ ). After the first round, we choose  $[z_1]$  as the smallest class in  $\mathcal{H}_1$ . Next, we destroy  $[z_1]$  by sealing  $z_1$ , which yields the augmented complex  $K'$ . This time, we choose  $[z_2]$ , giving  $\mathcal{H}_1 = \{[z_1], [z_2]\}$ .

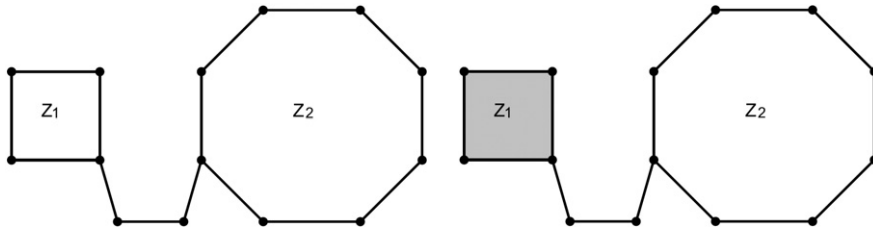


Fig. 4. Left: the original complex  $K$ . Right: the augmented complex  $K'$  after destroying the smallest class,  $[z_1]$ .

**Theorem 4.** *The procedure Measure-All( $K$ ) computes  $\mathcal{H}_d$ .*

**Proof.** We prove the theorem by showing that the sealing technique produces the same result as the naive greedy algorithm, namely,  $\mathcal{H}_d = \{h_{i_1}, h_{i_2}, \dots, h_{i_{\beta_d}}\}$ , assuming the classes are sorted according to size,  $S(h_{i_k}) \leq S(h_{i_{k+1}})$ . We show that for any  $l \in [0, \beta_d]$ , after computing and sealing the first  $l$  classes of  $\mathcal{H}_d$ , i.e.  $\{h_{i_1}, \dots, h_{i_l}\}$ , the next class we choose is exactly  $h_{i_{l+1}}$ . In other words, the localized cycle and size of the smallest class of the augmented complex  $K^l$  are equal to that of  $h_{i_{l+1}}$ .

First, any class between  $h_{i_l}$  and  $h_{i_{l+1}}$  in  $\text{seq}(H)$  will not be chosen. Any such class  $h_j$  is linearly dependent on classes that have already been chosen, namely,  $\{h_{i_1}, \dots, h_{i_l}\}$ . Since these classes have been sealed up, a cycle of  $h_j$  is a boundary in  $K^l$ . Thus,  $h_j$  cannot be chosen.

Second, according to algebra, one new cell can only destroy one class. Therefore, for any class in  $\text{seq}(H)$  that is not linearly dependent on  $\{h_{i_1}, \dots, h_{i_l}\}$ , it is nontrivial in  $K^l$ .

Third, the smallest class of  $K^l$ ,  $h_{\min}(K^l)$ , corresponds to  $h_{i_{l+1}}$ : any new simplex belonging to  $K^l \setminus K$  will not change the computation of the geodesic balls  $B_p^r$  with finite radius  $r$ , and thus will change neither the size measurement nor the localization. Thus, the  $h_{\min}(K^l)$  computed by the sealing technique is identical to  $h_{i_{l+1}}$  computed by the naive greedy method, in terms of the size and the localized cycle.  $\square$

#### 4.3. Complexity of the nonrefined algorithm

Throughout the algorithm, at most  $\beta_d$  new cells are added. The size of the augmented cell complex  $K'$  is  $O(n + \beta_d)$ . The procedure Measure-All( $K$ ) runs the procedure Measure-Smallest  $\beta_d$  times with  $K'$  as input. The procedure Measure-Smallest runs the persistent homology algorithm on  $K'$  using filter function  $f_p$  for each vertex of the original complex,  $K$ , and thus takes  $O(n(n + \beta_d)^3) = O(n^4)$  time. In total, it takes  $O(\beta_d n^4)$  time to compute the optimal basis.

### 5. An improvement using finite field linear algebra

In this section, we present an improvement on the algorithm presented in the previous section, more specifically, an improvement on computing the smallest geodesic ball carrying any nontrivial class (the procedure Measure-Smallest). The idea is based on the finite field linear algebra behind the homology.

In Section 5, we observe that for neighboring vertices,  $p_1$  and  $p_2$ , the birth times of the first essential homology class using  $f_{p_1}$  and  $f_{p_2}$  as filter functions are close (Theorem 6). This observation suggests that for each  $p$ , instead of computing  $B_p^{r(p)}$ , we may just test whether the geodesic ball centered at  $p$  with a certain radius carries any essential homology class. In Section 5.3, with some algebraic insight, we reduce the problem of testing whether a geodesic ball carries any essential homology class to the problem of comparing dimensions of two vector spaces. Furthermore, we use Lemma 7 to reduce the problem to rank computations of sparse matrices on the  $\mathbb{Z}_2$  field, for which we have ready tools from the literature. In Section 5.5, we conclude with detailed complexity analysis.

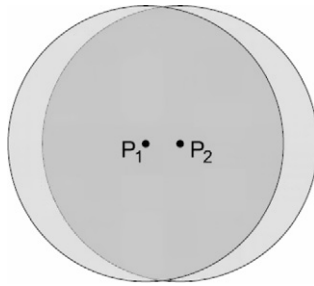
In this section, we will consider all edges to have weights of 1, for simplicity of exposition. However, please note that it is possible to generalize all results to deal with general (real) edge weights, though the algorithm becomes somewhat messier. We also assume that  $K$  has a single component; multiple components can be accommodated with a simple modification.

#### 5.1. Complexity

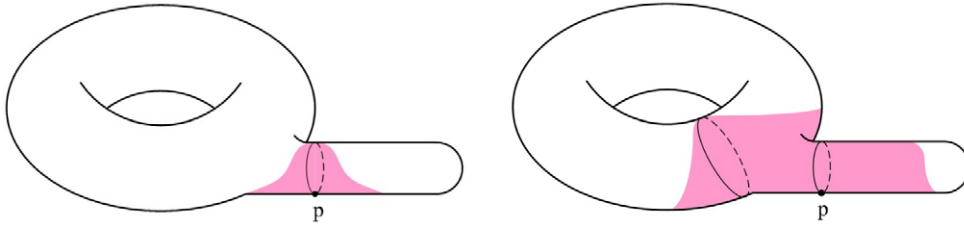
In doing so, we improve the complexity to  $O(\beta_d n^3 \log^2 n)$ . More detailed complexity analysis is provided in Section 5.5.

**Remark 5.** Cohen-Steiner et al. [7] provided a linear algorithm to maintain the persistences while changing the filter function. However, this algorithm is not directly applicable in our context. The reason is that it takes  $O(n)$  time to update the persistences for a transposition in the simplex-ordering. In our case, even for filter functions of two neighboring vertices, often it takes  $O(n^2)$  transpositions to transform one simplex-ordering into the other. See Fig. 5 for example. Therefore, updating the persistences while changing the filter function takes  $O(n^2) \times O(n) = O(n^3)$  time. This is the same amount of time it would take to compute the persistences from scratch.





**Fig. 5.** On the plane, when we change the filter function from  $f_{p_1}$  to  $f_{p_2}$ , in order to update the simplex-ordering, we should switch the order of the two blocks of simplices  $B_{p_1}^r \setminus B_{p_2}^r$  and  $B_{p_2}^r \setminus B_{p_1}^r$ , in which  $B_{p_1}^r$  and  $B_{p_2}^r$  are geodesic balls centered at  $p_1$  and  $p_2$  with a same radius,  $r$ . When  $r$  is big, these two blocks can have  $O(n)$  simplices. We then need  $O(n^2)$  transpositions to update the simplex-ordering.



**Fig. 6.** Only the ball in the second figure carries nonbounding cycles of  $K$ , although in both figures the balls have nontrivial topology.

5.2. Observation: neighboring vertices have similar geodesic distance functions

Since the filter functions of two neighboring vertices,  $f_{p_1}$  and  $f_{p_2}$ , are close to each other, the birth times of the first nonbounding cycles in both filters are close as well. This leads to Theorem 6. A simple proof is provided.

**Theorem 6.** *If two vertices  $p_1$  and  $p_2$  are neighbors, the birth times of the first nonbounding cycles for filter functions  $f_{p_1}$  and  $f_{p_2}$  differ by no more than 1.*

**Proof.**  $p_1$  and  $p_2$  are neighbors implies that for any point  $q$ ,

$$f_{p_2}(q) \leq f_{p_2}(p_1) + f_{p_1}(q) = 1 + f_{p_1}(q),$$

in which the inequality follows the triangular inequality. Therefore,  $B_{p_1}^{r(p_1)}$  is a subset of  $B_{p_2}^{r(p_1)+1}$ . The former carries nonbounding cycles implies that the latter does too, and thus  $r(p_2) \leq r(p_1) + 1$ . Similarly, we have  $r(p_1) \leq r(p_2) + 1$ .  $\square$

This theorem suggests a way to avoid computing  $B_p^{r(p)}$  for all  $p \in \text{vert}(K)$  in the procedure Measure-Smallest. Since our objective is to find the minimum of the  $r(p)$ 's, we do a breadth-first search through all the vertices with global variables  $r_{\min}$  and  $p_{\min}$  recording the smallest  $r(p)$  we have found and its corresponding center  $p$ , respectively. We start by applying the persistent homology algorithm on  $K$  with filter function  $f_{p_0}$ , where  $p_0$  is an arbitrary vertex of  $K$ . Initialize  $r_{\min}$  as the birth time of the first nonbounding cycle of  $K$ ,  $r(p_0)$ , and  $p_{\min}$  as  $p_0$ . Next, we do a breadth-first search through the rest vertices. For each vertex  $p_i, i \neq 0$ , there is a neighbor  $p_j$  we have visited (the parent vertex of  $p_i$  in the breath-first search tree). We know that  $r(p_j) \geq r_{\min}$  and  $r(p_i) \geq r(p_j) - 1$  (Theorem 6). Therefore,  $r(p_i) \geq r_{\min} - 1$ . We only need to test whether the geodesic ball  $B_{p_i}^{r_{\min}-1}$  carries any nonbounding cycle of  $K$ . If so,  $r_{\min}$  is decremented by one, and  $p_{\min}$  is updated to  $p_i$ . After all vertices are visited,  $p_{\min}$  and  $r_{\min}$  give us the ball we want.

However, testing whether the subcomplex  $B_{p_i}^{r_{\min}-1}$  carries any nonbounding cycle of  $K$  is not as easy as computing nonbounding cycles of the subcomplex. A nonbounding cycle of  $B_{p_i}^{r_{\min}-1}$  may not be nonbounding in  $K$  as we require. For example, in Fig. 6, the simplicial complex  $K$  is a torus with a tail. The shaded geodesic ball in the first figure does not carry any nonbounding cycle of  $K$ , although it carries its own nonbounding cycles. The geodesic ball in the second figure is the one that carries nonbounding cycles of  $K$ . Therefore, we need algebraic tools to distinguish nonbounding cycles of  $K$  from those of the subcomplex  $B_{p_i}^{r_{\min}-1}$ .

5.3. Procedure Contain-Nonbounding-Cycle: testing whether a subcomplex carries nonbounding cycles of  $K$

In this section, we present the procedure for testing whether a subcomplex  $K_0$  carries any nonbounding cycle of  $K$ . A chain in  $K_0$  is a cycle if and only if it is a cycle of  $K$ . However, solely from  $K_0$ , we are not able to tell whether a

cycle carried by  $K_0$  bounds or not in  $K$ . Instead, we write the set of cycles of  $K$  carried by  $K_0$ ,  $Z_d^{K_0}(K)$ , and the set of boundaries of  $K$  carried by  $K_0$ ,  $B_d^{K_0}(K)$ , as sets of linear combinations with certain constraints. Consequently, we are able to test whether any cycle carried by  $K_0$  is nonbounding in  $K$  by comparing the dimensions of  $Z_d^{K_0}(K)$  and  $B_d^{K_0}(K)$ . Lemma 7 shows that these dimensions can be computed by rank computations of sparse matrices.

To some extent, the idea of this section is similar in spirit to [27]. However, note that the two works developed independently.<sup>1</sup>

5.3.1. Expressing  $Z_d^{K_0}(K)$  and  $B_d^{K_0}(K)$  as sets of linear combinations with certain constraints

The set of boundaries and the set of cycles of  $K$  carried by  $K_0$  are

$$B_d^{K_0}(K) = B_d(K) \cap C_d(K_0) \quad \text{and}$$

$$Z_d^{K_0}(K) = Z_d(K) \cap C_d(K_0),$$

respectively. They are both vector spaces and the former is a subspace of the latter. It is not hard to show that the subcomplex  $K_0$  carries nonbounding cycles of  $K$  if and only if the dimensions of these two vector spaces are different. We now express them as linear combinations with certain constraints such that we can compute their dimensions using algebraic tools.

Let  $\hat{H}_d = [z_1, \dots, z_{\beta_d}]$  be the matrix whose column vectors are arbitrary  $\beta_d$  nonbounding cycles of  $K$  representing a homology basis. The boundary group and the cycle group of  $K$  are column spaces of the matrices  $\partial_{d+1}$  and  $\hat{Z}_d = [\partial_{d+1}, \hat{H}_d]$ , respectively.

$C_d(K_0)$  corresponds to the set of vectors each of whose  $i$ -th entry is zero for any simplex  $\sigma_i \notin K_0$ . We write  $Z_d^{K_0}(K)$  and  $B_d^{K_0}(K)$  as elements of  $Z_d(K)$  and  $B_d(K)$  whose  $i$ -th entries are zero. Consequently, we can write them as linear combinations with certain constraints,

$$B_d^{K_0}(K) = \{ \partial_{d+1} \gamma \mid \gamma \in \mathbb{Z}_2^{n(d+1)}, \partial_{d+1}^i \gamma = 0 \ \forall \sigma_i \notin K_0 \},$$

$$Z_d^{K_0}(K) = \{ \hat{Z}_d \gamma \mid \gamma \in \mathbb{Z}_2^{\beta_d + n(d+1)}, \hat{Z}_d^i \gamma = 0 \ \forall \sigma_i \notin K_0 \},$$

where  $\partial_{d+1}^i$  and  $\hat{Z}_d^i$  are the  $i$ -th rows of the matrices  $\partial_{d+1}$  and  $\hat{Z}_d$ , respectively. Here  $n(d+1)$  is the number of  $(d+1)$ -simplices in  $K$ , and thus the number of columns of  $\partial_{d+1}$ .

5.3.2. Computing dimensions by computing ranks of sparse matrices

With the following lemma, we can compute the dimensions of these two vector spaces  $Z_d^{K_0}(K)$  and  $B_d^{K_0}(K)$  by matrix rank computations. The proof is based on finite field linear algebra.

**Lemma 7.** For any matrix  $A = [A_1 A_2]$ ,  $\dim(\{A\gamma \mid A_2\gamma = 0\}) = \text{rank}(A) - \text{rank}(A_2)$ .

**Proof.** Denote  $P = \text{span } A = \{A\gamma\}$ .  $P_1 = \{A\gamma \mid A_2\gamma = 0\}$  is its subspace. The quotient vector space  $P/P_1$  is isomorphic to  $P_2 = \text{span}(A_2) = \{A_2\gamma\}$ . Therefore, we have

$$\begin{aligned} \dim(P_1) &= \dim(P) - \dim(P/P_1) \\ &= \dim(P) - \dim(P_2) \\ &= \text{rank}(A) - \text{rank}(A_2). \quad \square \end{aligned}$$

It is trivial to see that the order of the rows in these matrices does not interfere with the correctness of the theorem. The matrix  $A_2$  can be a certain subset of the rows of  $A$ , not necessarily the last few rows. Therefore, we can compute the dimensions of  $B_d^{K_0}(K)$  and  $Z_d^{K_0}(K)$  as

$$\dim(B_d^{K_0}(K)) = \text{rank}(\partial_{d+1}) - \text{rank}(\partial_{d+1}^{K \setminus K_0}), \quad \text{and}$$

$$\dim(Z_d^{K_0}(K)) = \text{rank}(\hat{Z}_d) - \text{rank}(\hat{Z}_d^{K \setminus K_0}),$$

where  $\partial_{d+1}^{K \setminus K_0}$  and  $\hat{Z}_d^{K \setminus K_0}$  are the matrices formed by rows of  $\partial_{d+1}$  and  $\hat{Z}_d$  whose corresponding simplices do not belong to  $K_0$ .

<sup>1</sup> The first draft of this paper was finished in April 2007.

### 5.4. Algorithm

The procedure Contain-Nonbounding-Cycle tests whether  $K_0$  carries any nonbounding cycle of  $K$  by testing whether  $\dim(B_d^{K_0}(K))$  and  $\dim(Z_d^{K_0}(K))$  are different. Since columns in  $\hat{H}_d$  correspond to  $\beta_d$  nonbounding cycles representing a homology basis, the ranks of  $\hat{Z}_d$  and  $\partial_{d+1}$  differ by  $\beta_d$ .  $K_0$  carries nonbounding cycles of  $K$  if and only if

$$\text{rank}(\hat{Z}_d^{K \setminus K_0}) - \text{rank}(\partial_{d+1}^{K \setminus K_0}) \neq \beta_d.$$

We use the algorithm of Wiedemann [22] for the rank computation.

In our algorithm, the boundary matrix  $\partial_{d+1}$  is given. The matrix  $\hat{H}_d$  can be precomputed by running persistent homology algorithm once, with an arbitrary filter function. Columns of  $\hat{H}_d$  are simply columns of matrix  $V$  corresponding to positive simplices paired with  $+\infty$ .

### 5.5. Complexity of the improved algorithm

The algorithm Measure-All( $K$ ) runs the improved procedure Measure-Smallest  $\beta_d$  times, with the augmented complex  $K'$  as the input complex. Measure-Smallest( $K'$ ) applies the persistent homology algorithm on  $K'$  once to compute  $\hat{H}_d$  and  $r(p_0)$ . Next, for each vertex, it runs the rank computation on submatrices of  $\partial_{d+1}$  and  $\hat{Z}_d = [\partial_{d+1}, \hat{H}_d]$ . Denoting  $m$  as the time of two rank computations, the algorithm takes  $O(\beta_d(n^3 + nm))$ , as the size of  $K'$  is  $O(n + \beta_d) = O(n)$ .

To know  $m$ , we need the number of nonzero entries in matrices  $\partial_{d+1}$  and  $\hat{Z}_d$ , as we are using a sparse matrix rank computation algorithm. Recall that in the augmented complex  $K'$ , we added  $O(\beta_d)$  new  $(d + 1)$ -dimensional cells, each of which has  $O(n)$   $d$ -faces. Therefore,  $\partial_{d+1}$  has  $O(n + \beta_d) = O(n)$  columns and  $O(n(d + 2) + n\beta_d) = O(nd + n\beta_d)$  nonzero entries. Since  $\hat{H}_d$  has  $\beta_d$  columns and  $O(n\beta_d)$  nonzero entries, the size and number of nonzero entries of  $\hat{Z}_d$  are asymptotically the same as  $\partial_{d+1}$ .

Running Wiedemann's rank computation on such matrices takes  $m = O(n \log n(nd + n\beta_d + n \log n))$ . If  $d$  and  $\beta_d$  are small enough – that is,  $O(\log n)$  or less – then we have improved the Measure-All( $K$ ) to  $O(\beta_d(n^3 + n^3 \log^2 n)) = O(\beta_d n^3 \log^2 n)$ . If we are dealing with an unusual situation in which  $d$  or  $\beta_d$  is large – say  $\Theta(n)$  – then the matrices are not sparse. We would prefer to use the old algorithm, with complexity  $O(\beta_d n^4)$ .

## 6. Stability result

In this section, we prove that our measurement of homology is stable: small changes of the geometry of the space imply small changes of our measurement. We define a change of the geometry of the space as a change of the metric in the space. We measure this change by measuring the  $L_\infty$ -norm difference of geodesic distance functions before and after the change. To facilitate the proof, we assume that during the change, the simplicial complex remains the same except in terms of its edge weights, and thus, the discrete geodesic distances. Formally, we quantify the change of the geometry as

$$\epsilon = \max_{p \in \text{vert}(K)} |f_p^1 - f_p^2|_\infty, \tag{1}$$

where  $f_p^1$  and  $f_p^2$  are the discrete geodesic distance functions from  $p$  before and after the change.

In this section, we prove the stability of our measurement by showing that (1) for a single homology class, the size is stable; and (2) for the whole homology group, although the optimal homology basis is not stable, the group structure filtered by the size is stable. For convenience, we drop the dimension of the pertinent homology,  $d$ , in notations.

### 6.1. A single class

For a single homology class, the size measure remains stable. Denote  $S^1(h)$  and  $S^2(h)$  as the size of class  $h$  before and after the change (computed using  $f_p^1$  and  $f_p^2$ , respectively). We have the following theorem.

**Theorem 8.**  $|S^1(h) - S^2(h)| \leq \epsilon$ , where  $\epsilon$  is the upper bound of the geometry change as defined in Eq. (1).

**Proof.** Denote  $r^1(p)$  and  $r^2(p)$  as the radii of the smallest geodesic balls carrying  $h$  computed using the geodesic distance  $f_p^1$  and  $f_p^2$ , respectively. We show that for any specific vertex  $p$ ,  $|r^1(p) - r^2(p)| \leq \epsilon$ . This leads to the fact that  $S^1(h) = \min_{p \in \text{vert}(K)} r^1(p)$  and  $S^2(h) = \min_{p \in \text{vert}(K)} r^2(p)$  differ in no more than  $\epsilon$ .

For any simplex  $\sigma$  in the ball  $B_p^{r^1(p)}$  calculated using  $f_p^1$ ,  $f_p^1(\sigma) \leq r^1(p)$ , and thus  $f_p^2(\sigma) \leq f_p^1(\sigma) + \epsilon \leq r^1(p) + \epsilon$ . This means that the ball  $B_p^{r^1(p)}$  calculated using  $f_p^1$  is a subcomplex of the ball  $B_p^{r^1(p)+\epsilon}$  calculated using  $f_p^2$ , which thus carries  $h$ . Therefore, according to the definition of  $r^2(p)$ , it is no greater than  $r^1(p) + \epsilon$ . Similarly,  $r^1(p) \leq r^2(p) + \epsilon$ .  $\square$

6.2. The homology group

Since the size of different classes can be very close, the optimal homology basis is not stable. For example, in Fig. 1 (Right), either  $\{[z_1], [z_2], [z_3]\}$  or  $\{[z_1], [z_2], [z_3] + [z_1]\}$  can be the optimal homology basis for little geometry changing, because the sizes of  $[z_3]$  and  $[z_1] + [z_3]$  are quite close. However, there is still some stability property in the homology group structure if we filter it with the class size. More specifically, the subgroup generated by small homology classes remains stable. For example, in Fig. 1 (Right), although the optimal homology basis is unstable, the subgroup generated by the two smaller classes in the optimal homology basis will always be the one generated by  $[z_1]$  and  $[z_2]$ .

We formalize this stability by defining the subgroup filtration of a topological space and the distance between two such filtrations. A subgroup filtration is a sequence of subgroups of the homology group generated by subsets of the optimal homology basis filtered by the class size. A formal definition is as follows.

**Definition 9 (Subgroup filtration).** Given an optimal homology basis  $\mathcal{H} = \{h_1, h_2, \dots, h_\beta\}$ , where we assume  $S(h_i) \leq S(h_{i+1})$ , a *subgroup filtration* is a sequence of subgroups of the homology group,  $\mathcal{X} = \{\psi_0, \psi_1, \psi_2, \dots, \psi_\beta\}$ , where  $\psi_i = \text{span}(h_1, h_2, \dots, h_i)$  is the subgroup generated by the classes  $h_1, h_2, \dots, h_i$ .

Since here the homology group and all its subgroups are vector spaces, we use the notation  $\psi_i = \text{span}(h_1, h_2, \dots, h_i)$  when we say  $h_1, h_2, \dots, h_i$  generates  $\psi_i$ .

Obviously, the subgroup filtration is a sequence of subgroups of  $H(K)$  with a nested structure

$$\emptyset = \psi_0 \subset \psi_1 \subset \dots \subset \psi_\beta = H(K).$$

For convenience, we denote the size of a subgroup,  $\psi_i$ , as the size of the largest class in the optimal homology basis generating  $\psi_i$ , formally,  $S(\psi_i) = S(h_i)$ . Please note that  $S(\psi_i)$  is not the size of the largest class in  $\psi_i$ .

Given two different sets of discrete geodesic distance functions (different geometries) of a same topological space,  $f^1$  and  $f^2$ , we have two different subgroup filtrations  $\mathcal{X}^1$  and  $\mathcal{X}^2$ . Next, we define their distance, which requires the definition of the projection of one subgroup in one filtration onto the other filtration.

**Definition 10 (Projection).** Given two subgroup filtrations of a same homology group  $\mathcal{X}^1 = \{\psi_0^1, \psi_1^1, \psi_2^1, \dots, \psi_\beta^1\}$  and  $\mathcal{X}^2 = \{\psi_0^2, \psi_1^2, \psi_2^2, \dots, \psi_\beta^2\}$ , define the *projection* of  $\psi_i^1$  onto  $\mathcal{X}^2$  as the first subgroup in  $\mathcal{X}^2$  that carries  $\psi_i^1$ , formally,

$$\text{proj}(\psi_i^1, \mathcal{X}^2) = \psi_j^2 \quad \text{s.t.} \quad j = \min_{\psi_i^1 \subseteq \psi_k^2} k.$$

**Definition 11 (Distance).** Define the *distance* between  $\mathcal{X}^1$  and  $\mathcal{X}^2$  as the maximal difference between the sizes of any subgroup in  $\mathcal{X}^1$  or  $\mathcal{X}^2$  and its projection onto the other filtration, formally,

$$\text{dist}(\mathcal{X}^1, \mathcal{X}^2) = \max \left\{ \max_i |S^1(\psi_i^1) - S^2(\text{proj}(\psi_i^1, \mathcal{X}^2))|, \max_i |S^2(\psi_i^2) - S^1(\text{proj}(\psi_i^2, \mathcal{X}^1))| \right\}.$$

Let  $\mathcal{X}^1$  and  $\mathcal{X}^2$  be the subgroup filtrations of the original space and the one after the change. We can prove the following stability result.

**Theorem 12.**

$$\text{dist}(\mathcal{X}^1, \mathcal{X}^2) \leq \epsilon = \max_{p \in \text{vert}(K)} |f_p^1 - f_p^2|_\infty. \tag{2}$$

**Proof.** Take a subgroup  $\psi_i^1$ , generated by  $h_1^1, h_2^1, \dots, h_i^1$ , the smallest  $i$  elements of the optimal homology basis  $\mathcal{H}^1$ , determined by  $f^1$ . For any  $j \in [1, i]$ , we have

$$S^2(h_j) \leq S^1(h_j) + \epsilon \leq S^1(\psi_i^1) + \epsilon,$$

in which the first and the second inequalities are due to Theorem 8 and the definition of  $\psi_i^1$ , respectively. Therefore, we have

$$S^2(\text{proj}(\psi_i^1, \mathcal{X}^2)) \leq \max_{j \in [1, i]} S^2(h_j) \leq S^1(\psi_i^1) + \epsilon.$$

This is true for all  $i \in [1, \beta]$ . Similarly we can prove for any subgroup of  $\mathcal{X}^2$ , its distance from its projection onto  $\mathcal{X}^1$  is upper-bounded by  $\epsilon$ . Eq. (2) is proved.  $\square$

## 7. Conclusion

In this paper, we defined a size measure for homology classes. We provided an algorithm to compute the optimal homology basis, namely, the basis whose elements have the minimal total size. Using finite field linear algebra, we improved the complexity of our straightforward algorithm. Finally, we proved that our size measure is stable in a natural sense.

### 7.1. Measure a given class

One interesting question is, instead of computing the optimal basis, can we measure a single given class,  $[z]$ . We modify the procedure Measure-Smallest to achieve this. Again, we iterate through all vertices. For each vertex  $p$ , we find the smallest geodesic ball centered at  $p$  carrying  $[z]$ , namely,  $B_p^{r(p)}$ . We apply persistent homology on the complex using  $f_p$  as the filter function. We pick all the columns in  $V$  corresponding to positive simplices that are paired with  $+\infty$ , namely,  $z_1, z_2, \dots, z_{\beta_d}$ , sorted according to their order in the filtration. We find the smallest index  $i$  so that  $z$  is a linear combination of boundaries and  $z_1, z_2, \dots, z_i$ , namely,

$$z = [\partial_{d+1}, z_1, z_2, \dots, z_i]\gamma. \quad (3)$$

The positive simplex corresponding to this smallest  $i$  gives us  $r(p)$ . Replacing  $\partial_{d+1}$  with 0, we get a representative cycle of  $[z]$  carried by  $B_p^{r(p)}$ ,  $[0, z_1, z_2, \dots, z_i]\gamma$ . Iterating through every vertex  $p$ , we find the smallest ball carrying  $[z]$ ,  $B_{p_{\min}}^{r(p_{\min})}$ , and consequently the size and localized cycle of  $[z]$ .

## Acknowledgements

We appreciate constructive comments from anonymous reviewers.

## References

- [1] Gunnar Carlsson, Persistent homology and the analysis of high dimensional data, in: Fields-Ottawa Workshop on the Geometry of Very Large Data Sets, February 2005.
- [2] Robert Ghrist, Barcodes: The persistent topology of data, *Bull. Amer. Math. Soc.* 45 (1) (2008) 61–75.
- [3] Jeff Erickson, Sarel Har-Peled, Optimally cutting a surface into a disk, *Discrete Comput. Geom.* 31 (1) (2004) 37–59.
- [4] Christopher Carner, Miao Jin, Xianfeng Gu, Hong Qin, Topology-driven surface mappings with robust feature alignment, in: Proceedings of the 16th IEEE Visualization Conference, 2005, pp. 543–550.
- [5] Vin de Silva, Robert Ghrist, Coverage in sensor networks via persistent homology, *Algebr. Geom. Topol.* 7 (2007) 339–358.
- [6] Pankaj K. Agarwal, Herbert Edelsbrunner, John Harer, Yusu Wang, Extreme elevation on a 2-manifold, *Discrete Comput. Geom.* 36 (2006) 553–572.
- [7] David Cohen-Steiner, Herbert Edelsbrunner, Dmitriy Morozov, Vines and vineyards by updating persistence in linear time, in: Proceedings of the 22nd ACM Symposium on Computational Geometry, 2006, pp. 119–126.
- [8] Zoë J. Wood, Hugues Hoppe, Mathieu Desbrun, Peter Schröder, Removing excess topology from isosurfaces, *ACM Trans. Graph.* 23 (2) (2004) 190–208.
- [9] Gunnar Carlsson, Tigran Ishkhanov, Vin de Silva, Leonidas J. Guibas, Persistence barcodes for shapes, *Internat. J. Shape Modeling* 11 (2) (2005) 149–188.
- [10] Herbert Edelsbrunner, David Letscher, Afra Zomorodian, Topological persistence and simplification, *Discrete Comput. Geom.* 28 (4) (2002) 511–533.
- [11] Afra Zomorodian, *Topology for Computing*, Cambridge Monogr. Appl. Comput. Math., Cambridge University Press, 2005.
- [12] Afra Zomorodian, Gunnar Carlsson, Computing persistent homology, *Discrete Comput. Geom.* 33 (2) (2005) 249–274.
- [13] Gunnar Carlsson, Afra Zomorodian, The theory of multidimensional persistence, in: Proceedings of the 23rd ACM Symposium on Computational Geometry, 2007, pp. 184–193.
- [14] David Cohen-Steiner, Herbert Edelsbrunner, John Harer, Stability of persistence diagrams, *Discrete Comput. Geom.* 37 (2007) 103–120.
- [15] Herbert Edelsbrunner, John Harer, Persistent homology – A survey, in: J.E. Goodman, J. Pach, R. Pollack (Eds.), *Twenty Years After*, AMS, 2007.
- [16] David Cohen-Steiner, Herbert Edelsbrunner, John Harer, Extending persistent homology using Poincaré and Lefschetz duality, *Found. Comput. Math.* 9 (1) (2009) 79–103.
- [17] Afra Zomorodian, Gunnar Carlsson, Localized homology, in: Proceedings of the 2007 International Conference on Shape Modeling and Applications, 2007, pp. 189–198.
- [18] Jeff Erickson, Kim Whittlesey, Greedy optimal homotopy and homology generators, in: Proceedings of the 16th Annual ACM-SIAM Symposium on Discrete Algorithms, 2005, pp. 1038–1046.
- [19] Igor Guskov, Zoë J. Wood, Topological noise removal, in: Proceedings of the Graphics Interface 2001 Conference, 2001, pp. 19–26.
- [20] J.R. Munkres, *Elements of Algebraic Topology*, Addison-Wesley, Redwood City, CA, 1984.
- [21] T. Kaczynski, M. Mrozek, M. Slusarek, Homology computation by reduction of chain complexes, *Comput. Math. Appl.* 35 (1998) 59–70.
- [22] Douglas H. Wiedemann, Solving sparse linear equations over finite fields, *IEEE Trans. Inform. Theory* 32 (1) (1986) 54–62.
- [23] Chao Chen, Daniel Freedman, Quantifying homology classes, in: Proceedings of the 25th Annual Symposium on Theoretical Aspects of Computer Science, 2008, pp. 169–180.
- [24] Chao Chen, Daniel Freedman, Hardness results for homology localization, preprint, June 2009.
- [25] T.H. Cormen, C.E. Leiserson, R.L. Rivest, C. Stein, *Introduction to Algorithms*, MIT Press, 2001.
- [26] James G. Oxley, *Matroid Theory*, Oxford University Press, 1992.
- [27] David Cohen-Steiner, Herbert Edelsbrunner, John Harer, Dmitriy Morozov, Persistent homology for kernels, images, and cokernels, in: *SODA*, 2009, pp. 1011–1020.