

# 1 APPENDIX

## 2 Homology of Simplicial Complexes

3 Listed here are the basic concepts in algebraic topology which are necessary in understanding of per-  
4 sistent homology. Definitions and theorems are taken mainly from Bubenik (2015), Carlsson (2009),  
5 Edelsbrunner and Harer (2008), Ghrist (2008), Pun et al. (2018) Otter et al. (2017) and Zomorodian and  
6 Carlsson (2005).

### 7 **Simplices.**

8 One way of associating an algebraic and combinatorial structure to a topological space is by use of  
9 simplicial complexes.

10 **Definition 1** A *k-simplex*,  $\Delta^k$ , is the convex hull of  $k + 1$  points which do not lie in a hyperplane of  
11 dimension  $k$  or less. It can also be denoted as  $[v_0, v_1, v_2, \dots, v_k]$ , where  $v_i$ 's are the vertices of  $\Delta^k$  which  
12 has the natural ordering and  $k$  is its dimension.

13 **Definition 2** A *face* of a  $k$ -simplex  $[v_0, v_1, v_2, \dots, v_k]$  is a simplex  $[v_{i_1}, v_{i_2}, \dots, v_{i_k}]$  where  $i_j \in \{0, 1, 2, \dots, k\}$   
14 for each  $j$  and  $0 \leq i_1 < i_2 < i_3 < \dots < i_k \leq k$ . If a simplex  $\sigma'$  is a face of a simplex  $\sigma$ , then it is denoted  
15 as  $\sigma' \subseteq \sigma$ . If the dimension of  $\sigma'$  is less than the dimension of  $\sigma$ , then  $\sigma' \subset \sigma$ .

16 A 0-simplex can be represented as a point, a 1-simplex as an edge from one vertex to another vertex, a  
17 2-simplex as a triangular region defined by 3 non-collinear points, a 3-simplex as a tetrahedron together  
18 with its interior defined by 4 non-coplanar points, and so on.

### 19 **Simplicial Complexes.**

20 The  $k$ -simplices are regarded as building blocks of simplicial complexes. Simplices can be glued together  
21 to form simplicial complexes. A simplicial complex is formally defined as follows.

22 **Definition 3** Let  $K^0$  be a set of vertices. A **simplicial complex**,  $K$ , is a collection of simplices whose  
23 vertices are element of  $K^0$ , such that if  $v \in K^0$  then  $[v] \in K$  and if  $\tau$  and  $\sigma$  are simplices such that  $\sigma \in K$   
24 and  $\tau \subset \sigma$  then  $\tau \in K$ . Moreover, the dimension of  $K$  is the maximum of the dimensions of its elements.

25 **Definition 4** Let  $K^i$  denote the set of  $i$ -simplices in a simplicial complex  $K$ . The *n-skeleton* of  $K$  is the  
26 union of the sets  $K^i$  for all  $i \in \{0, 1, 2, \dots, n\}$ . If  $\sigma_1$  is a simplex of dimension  $n_1$  and  $\sigma_2$  is a simplex of  
27 dimension  $n_2$ , such that  $\sigma_1 \subset \sigma_2$ , then  $\sigma_1$  is said to be a face of  $\sigma_2$  of codimension  $n_2 - n_1$ .

### Example 1 Simplicial Complex A

$$A = \{[a], [b], [c], [d], [e], [a, b], [b, c], [a, c], [d, a], [c, d], [a, b, c]\}$$

### Example 2 Simplicial Complex B

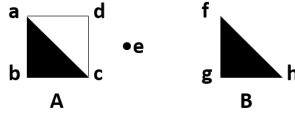
$$B = \{[f], [g], [h], [f, g], [g, h], [f, h], [f, g, h]\}$$

28 Simplicial complex  $A$  is of dimension 2 since it contains a 2-simplex and it is the element of  $A$  with  
29 the largest dimension. Similarly,  $B$  is of dimension 2.

30 A simplicial complex can be referred to as an abstract simplicial complex because of its abstract  
31 nature. But, one can interpret a finite simplicial complex geometrically as a subset of  $\mathbb{R}^n$  for some natural  
32 number  $n$ . Such subset is called a geometric realization and it is unique up to a canonical piecewise-linear  
33 homomorphism (Otter et al., 2017). That is, for a simplicial complex  $K$ , there exists a geometric simplicial  
34 complex  $G$  whose vertices are in one-to-one correspondence with the vertices of  $K$  and a subset of vertices  
35 in  $K$  define a simplex in  $G$  if and only if they correspond to the vertices of some simplex of  $K$ .

36 Figure 1 shows the respective geometric realization of simplicial complexes  $A$  and  $B$  in  $\mathbb{R}^2$ .

37 Note that a simplicial complex  $\Delta$  can also be viewed as a topological space expressed as a quotient of  
38 disjoint union of simplices by an equivalence relation that identifies certain faces of certain simplices.



**Figure 1.** Geometric Realization of A and B

39 **Homology of Simplicial Complexes.**

40 A formal sum of  $k$ -simplices is called a  $k$ -chain and the free abelian group having a collection of  
 41  $k$ -simplices as its basis is called a chain group.

Let  $X$  be a simplicial complex and  $\Delta_k(X)$  be the free abelian group generated by the  $k$ -simplices of  $X$ . Elements of  $\Delta_k(X)$  are called  $k$ -simplicial chains. For any  $k \in \{1, 2, 3, \dots\}$ , define the boundary map as the linear map

$$\partial_k : \Delta_k(X) \rightarrow \Delta_{k-1}(X),$$

$$\sigma \mapsto \sum_{\tau \subset \sigma, \tau \in \Delta^{n-1}} \tau.$$

42 The boundary map  $\partial_k$  maps each  $k$ -simplex to its boundary, which is the sum of its faces of codimension  
 43 1. The map  $\partial_0$  is called the zero map. It can be shown that  $\partial_n \circ \partial_{n+1} = 0$ , that is the boundary of a boundary  
 44 is always empty. Moreover, the image of  $\partial_{n+1}$  is contained in the kernel of  $\partial_k$ .

The boundary operators and the chain groups form into a chain complex  $\mathbf{C}_*$ :

$$\dots \rightarrow \Delta_{k+1} \xrightarrow{\partial_{k+1}} \Delta_k \xrightarrow{\partial_k} \Delta_{k-1} \rightarrow \dots$$

**Definition 5** For each  $n \in \{0, 1, 2, 3, \dots\}$ , the  $n$ -th homology of a simplicial complex  $X$ , is given as

$$H_n(X) := \text{Ker}(\partial_n) / \text{Im}(\partial_{n+1}).$$

Moreover, its dimension

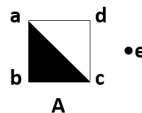
$$\beta_n(X) := \dim H_n(X) = \dim \text{Ker}(\partial_n) - \dim \text{Im}(\partial_{n+1})$$

45 is called the  $n$ -th Betti number of  $X$ , or the rank of the  $n$ -th homology group of  $(X)$ . And, elements of  
 46  $\text{Im}(\partial_{n+1})$  are called  $n$ -boundaries, and elements of  $\text{Ker}(\partial_n)$  are called  $n$ -cycles.

The  $n$ -cycles which are not boundaries represent  $n$ -dimensional holes. Thus, the  $n$ -th Betti number gives the number of  $n$ -holes. Particularly, the  $\beta_0(X)$  gives the number of connected components, the  $\beta_1(X)$  gives the number of tunnels, the  $\beta_2(X)$  gives the number of voids, and so on. Furthermore, if  $X$  is a simplicial complex of dimension  $p$ , then  $H_n(X) = 0$  for each  $n > p$ . Then there is the following sequence,

$$0 \xrightarrow{\partial_{n+1}} \Delta_n(X) \xrightarrow{\partial_n} \dots \xrightarrow{\partial_2} \Delta_1(X) \xrightarrow{\partial_1} \Delta_0(X) \xrightarrow{\partial_0} 0.$$

47 **Example 3** Consider the simplicial complex  $A = \{[a], [b], [c], [d], [e], [a, b], [b, c], [a, c], [d, a], [c, d], [a, b, c]\}$   
 48 from Example 1 with the geometric realization given in Fig. 2.



**Figure 2.** Geometric Realization of A and B

Then there is the following sequence,

$$0 \dots 0 \xrightarrow{\partial_3} \Delta_2(A) \xrightarrow{\partial_2} \Delta_1(A) \xrightarrow{\partial_1} \Delta_0(A) \xrightarrow{\partial_0} 0,$$

49 where

$$\begin{aligned} 50 \Delta_0(A) &= \mathbb{Z}^5 = \text{span}_{\mathbb{Z}}\{[a], [b], [c], [d], [e]\}, \\ 51 \Delta_1(A) &= \mathbb{Z}^5 = \text{span}_{\mathbb{Z}}\{[a, b], [b, c], [a, c], [c, d], [d, a]\}, \\ 52 \Delta_2(A) &= \mathbb{Z} = \text{span}_{\mathbb{Z}}\{[a, b, c]\}, \text{ and} \\ 53 \Delta_k(A) &= 0 \text{ for each } k \geq 3. \end{aligned}$$

54  
55 Also, for  $k = 1, 2, 3$ , the boundary operator  $\partial_k$  is defined for  $k$ -simplices, respectively as follows,  
56  $\partial_0([x]) = 0$  for each  $[x] \in \Delta_0(A)$ ,  
57  $\partial_1([x, y]) = [y] - [x]$  for each  $[x, y] \in \Delta_1(A)$ , and  
58  $\partial_2([x, y, z]) = [x, y] + [y, z] - [x, z]$  for each  $[x, y, z] \in \Delta_2(A)$ .

59  
60 The homology groups are computed as follows,

$$\begin{aligned} 61 H_0(A) &= \frac{\ker \partial_0}{\text{im} \partial_1} \\ 62 &= \frac{\text{span}_{\mathbb{Z}}\{[a], [b], [c], [d], [e]\}}{\text{span}_{\mathbb{Z}}\{([b] - [a]) + ([c] - [b]), ([c] - [a]) + ([d] - [c]), [a] - [d] + [c] - [a]\}} \\ 63 &= \frac{\text{span}_{\mathbb{Z}}\{[a], [b], [c], [d], [e]\}}{\text{span}_{\mathbb{Z}}\{[a] - [c], [a] - [d], [d] - [c]\}} = \frac{\mathbb{Z}^5}{\mathbb{Z}^3} = \mathbb{Z}^2, \\ 64 H_1(A) &= \frac{\ker \partial_1}{\text{im} \partial_2} = \frac{\text{span}_{\mathbb{Z}}\{[c, d], [d, a]\}}{\text{span}_{\mathbb{Z}}\{[a] + [b] - [c]\}} = \frac{\mathbb{Z}^2}{\mathbb{Z}^1} = \mathbb{Z}, \\ 65 H_2(A) &= \frac{\ker \partial_2}{\text{im} \partial_3} = 0, \text{ and} \\ 66 H_k(A) &= 0 \text{ for each } k \geq 3. \end{aligned}$$

67  
68 The Betti numbers are  $\beta_0 = 2$ ,  $\beta_1 = 1$ ,  $\beta_2 = 0$ , which means that there are 2 connected spaces, 1 hole  
69 and 0 voids in  $A$ .

70 For the succeeding sections, simplicial homology will be defined over the field  $\mathbb{F}_2$  with 2 elements,  
71 where  $1 \neq -1$ . So instead of defining the chain groups as free abelian groups, we define the chain groups  
72 as vector spaces over  $\mathbb{F}_2$ . However, when computing simplicial homology over  $\mathbb{F}_2$ , one needs to be  
73 careful when defining the boundary maps  $\partial_k$  to ensure that  $\partial_k \circ \partial_{k+1}$  remains the zero map (Otter et al.,  
74 2017). Consequently, the definition is just almost the same, but the resulting homology groups and Betti  
75 numbers may vary for different fields. For the purpose of using persistent homology in data science or on  
76 Euclidean spaces, it suffices to consider homology with coefficients in the field  $\mathbb{F}_2$ . Indeed, we will see in  
77 the discussion of obtaining topological summaries in the form of barcodes that we will need to compute  
78 homology with coefficients in a field, particularly  $\mathbb{F}_2$ . Furthermore, most of the implementations for the  
79 computation of persistent homology in the examples work with  $\mathbb{F}_2$ .

## 80 **Computing Persistent Homology of a Point Cloud**

81 Presented here is the general guideline in computing persistent homology of a dataset which follows the  
82 pipeline of computing persistent homology as presented in Otter et al. (2017).

83 Data can be viewed as a collection of points in a metric space. This finite metric space is also called a  
84 point cloud. Points in the dataset or point cloud are thickened gradually and this gradual evolution of the  
85 point clouds' shape and its topological properties are now the point of interest for topologist and data  
86 scientists.

87 Let  $X_0$  be a finite subset of a Euclidean space  $R^n$ . Then,  $X_0$  is an example of a point cloud. These  
88 points which can be viewed as vertices which will serve as building blocks of complexes.

89 Let  $\epsilon$  be a non-negative real number which serves as parameter to thicken  $X_0$  and  $X_\epsilon$  be the thickened  
90 point cloud. Homology of a given data pertains to topological invariant properties of  $X_\epsilon$  which can be  
91 computed algebraically. That is, for each nonnegative integer  $i$ , there is a corresponding vector space or a  
92 homology group  $H_i(X_\epsilon)$ . The dimension of the first 3 homology groups gives the number of connected  
93 components, the number of tunnels or holes and the numbers of voids, respectively. This algebraic  
94 structures are said to be robust or homotopy invariant. That is, a primary space's topological invariant  
95 features does not change when the space undergo bending, stretching or other deformations. Furthermore,  
96 computing homology of a finite simplicial complex can be easily done with the aid of linear algebra.  
97

98 Computing the homology of an arbitrary topological space is not as straight forward as that of  
 99 computing the homology of a finite simplicial complex. Firstly, one has to find a simplicial complex  
 100 whose homology approximates the homology of the arbitrary space. There are various ways of doing this,  
 101 and the most natural methods are with the use of Čech complexes and Vietoris-Rips complexes. There are  
 102 alternative complexes like Delaunay complex, Alpha complex and Witness complex.

103 Let the parameter  $\varepsilon$  be a non-negative real number and  $X$  be a set of points in the Euclidean space  $\mathbb{E}^k$   
 104 and  $\mathcal{U} = \{U_i\}_{i \in I}$  be a cover of  $X$ . The  $k$ -simplices of the Čech complex are the non-empty intersections  
 105 of  $k+1$  sets in  $\mathcal{U}$ .

106 **Definition 6** Let  $\mathcal{U} = \{U_i\}_{i \in I}$  be the non-empty collection of sets. The nerve of  $\mathcal{U}$  is the simpli-  
 107 cial complex with the vertices given by  $I$  and the  $k$ -simplices given by  $\{i_0, i_1, i_2, \dots, i_k\}$  if and only if  
 108  $\bigcap_{j \in \{0, 1, \dots, k\}} U_{i_j} \neq \emptyset$ .

109 **Theorem 1 (Nerve Theorem)** The geometric realization of the nerve of  $\mathcal{U}$  is homotopy equivalent to  
 110 the union of sets in  $\mathcal{U}$ .

**Definition 7** The Čech complex with parameter  $\varepsilon$  of  $X$  is given as

$$\check{C}_\varepsilon(X) := \left\{ \sigma \in X \mid \bigcap_{x \in \sigma} B(x, \varepsilon) \neq \emptyset \right\}$$

111 where  $B(x, \varepsilon)$  is a closed ball of radius  $\varepsilon$  centered at  $x$ .

112 If the cover of the sets in  $X$  is sufficiently ‘nice,’ then the Nerve Theorem guarantees that the nerve of  
 113 the cover and the space  $X$  have the same homology (Edelsbrunner and Harer, 2010). However, finding  
 114 the Čech complex is computationally expensive as it involves investigating a very large number of  
 115 intersections. Furthermore, the Čech complex may have a higher dimension than the underlying space.  
 116 This is the reason why choosing Vietoris-Rips (VR) complex, an approximation of the Čech complex, can  
 117 be more attractive.

**Definition 8** Let  $(X, d)$  be a metric space,  $S$  be a subspace of  $X$  with the induced metric. The Vietoris-Rips  
 complex with parameter  $\varepsilon$ , denoted by  $\mathcal{R}_\varepsilon(X)$ , is the set of all  $\sigma \subset X$ , such that the largest Euclidean  
 distance between any of its points is at most  $2\varepsilon$ . That is, given  $S \subset X$ ,

$$\mathcal{R}_\varepsilon(S) = \{ \sigma \subseteq S \mid d(x, y) \leq 2\varepsilon \text{ for all } x, y \in \sigma \}.$$

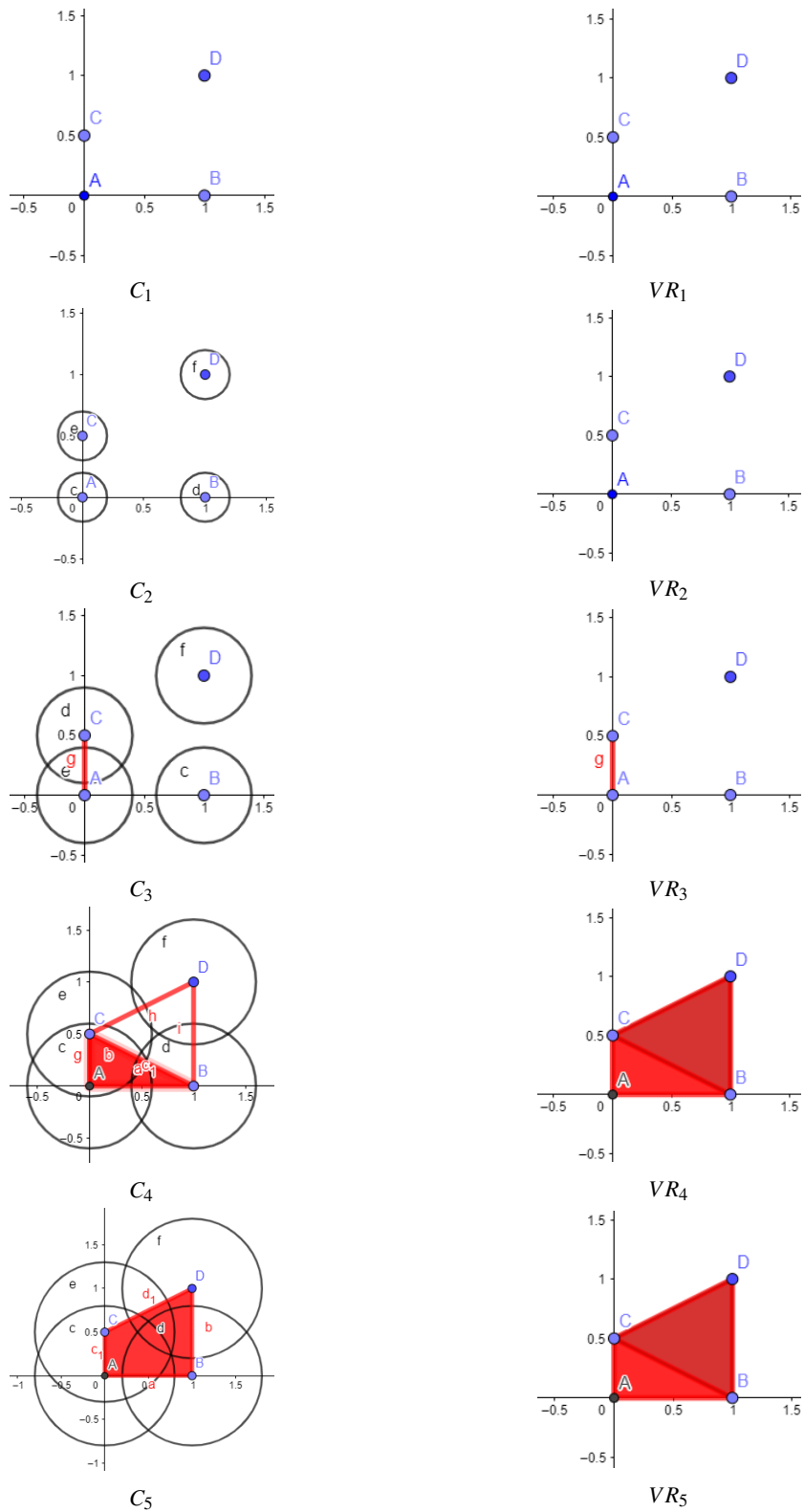
118 Both the Vietoris-Rips complex and the Čech complex are abstract simplicial complexes which may  
 119 be defined at various parameters  $\varepsilon$ , but only Čech complex preserves the homotopy information of the  
 120 topological spaces formed by the  $\varepsilon$ -balls.

121 Moreover, given  $S$  is a subset of a Euclidean space, the Vietoris-Rips complex approximates the Čech  
 122 complex in such a way that  $\check{C}_\varepsilon(S) \subseteq VR_\varepsilon(S) \subseteq \check{C}_{\sqrt{2}\varepsilon}(S)$ .

123 The construction of a VR complex can be made easier with the use of clique complexes, also known  
 124 as the flag complexes. In topology, recall that a graph is complete if any two vertices in the graph is  
 125 connected by an edge and the set of vertices which form a complete graph is called a clique. A  $k$ -clique  
 126 complex is formed from a clique of  $k+1$  vertices. Since subsets of a clique is also a clique, then a clique  
 127 complex is also a simplicial complex.

128 Now, for easier construction a VR complex of  $S$ , form the  $\varepsilon$ -neighborhood of  $S$  which is composed  
 129 of vertices in  $S$  and edges  $(i, j) \in S \times S$  where  $i \neq j$  and  $d(i, j) \leq 2\varepsilon$ . Afterwards, compute the clique  
 130 complex of the  $\varepsilon$ -neighborhood graph. This construction is easier for the reason that one only needs to  
 131 check for pairwise distances to construct a clique complex. Although this technique is computationally  
 132 less expensive as that of computing the Čech complex, the VR complex may have the same worst-case  
 133 complexity as that of the Čech complex. That is, in the worst case, a VR complex may have  $2^{|S|} - 1$   
 134 simplices and dimension  $|S| - 1$ . Moreover, one can opt to just compute the VR complex up to some  
 135 dimension  $k \ll |S| - 1$ . Otter et al. (2017) used  $k = 2$  and  $k = 3$  in their simulations.

136 **Example 4** Consider  $K = \{(0, 0), (1, 0), (0, 0.5), (1, 1)\} \subseteq \mathbb{R}^2$  and the filtration parameters values  $\varepsilon_1 = 0$ ,  
 137  $\varepsilon_2 = 0.2$ ,  $\varepsilon_3 = 0.4$ ,  $\varepsilon_4 = 0.6$  and  $\varepsilon_5 = 0.8$ . The first column of Fig. 3 shows the filtration of  $K$  using Čech  
 138 complexes, with  $C_1 \subseteq C_2 \subseteq \dots \subseteq C_5$ . The second column of Fig. 3 show the filtration of  $K$  using  
 139 Vietoris-Rips complexes, with  $VR_1 \subseteq VR_2 \subseteq \dots \subseteq VR_5$ .



**Figure 3.** Filtration of  $K$  using Čech and Vietoris-Rips Complexes

140 In this particular example, the Čech complexes and the Vietoris-Rips complexes differ only at the  
 141 4-th filtration index.

142 Given a finite metric space  $S$ , say a set of experimental dataset, it is assumed that the data is a sample  
 143 from some underlying topological space. The goal in computing the persistent homology of this data  
 144 is to recover the properties of such underlying topological space while maintaining the robustness of  
 145 the data. Given a subset  $S$  of a Euclidean space, one can consider  $S_\varepsilon$ , a simplicial complex at different  
 146 values of  $\varepsilon$ . As the value of  $\varepsilon$  increases, simplices are added to the complexes and sequence of nested  
 147 simplicial complexes is formed. This is the part where the homology of the simplicial complexes changes  
 148 as the parameter  $\varepsilon$  changes. The goal now in is to determine which topological features persist across the  
 149 changes in the  $\varepsilon$  values. Thus, the name persistent homology.

150 **Definition 9** Let  $K$  be a finite simplicial complex and  $K_1 \subseteq K_2 \subseteq \dots \subseteq K_r = K$  be a finite sequence of  
 151 nested subcomplexes of  $K$ .  $K$  is called a filtered simplicial complex and the sequence  $\{K_1, K_2, \dots\}$  is called  
 152 the filtration of  $K$ .

153 Homology of each of the subcomplexes can be computed. For each  $p$ , the inclusion maps  $K_i \rightarrow$   
 154  $K_j$  induce  $\mathbb{F}_2$ -linear maps  $\partial_i^j : H_p(K_i) \rightarrow H_p(K_j)$  for all  $i, j \in \{1, 2, \dots, r\}$  with  $i \leq j$ . It follows from  
 155 functoriality that  $\partial_k^j \circ \partial_i^k = \partial_i^j$  for all  $i \leq k \leq j$ .

156 **Definition 10** Let  $K_s$  be a subcomplex in the filtration of the simplicial complex  $K$ , or  $K_s$  be the filtered  
 157 complex at time  $s$ , and  $Z_k^s = \text{Ker} \partial_k^s$  and  $B_k^s = \text{Im} \partial_{k+1}^s$  be the  $k$ -th cycle group and boundary group of  $K_s$ ,  
 158 respectively. The  $k$ -th homology group of  $K_s$  is  $H_k^s = Z_k^s / B_k^s = \text{Ker}(\partial_k^s) / \text{Im}(\partial_{k+1}^s)$ .

**Definition 11** For  $p \in \{0, 1, 2, \dots\}$ , the  $p$ -persistent  $k$ -th homology group of  $K$  given a subcomplex  $K_s$  is  
 $K_s$  is

$$H_k^{s,p}(K, K_s) = H_k^{s,p}(K) = Z_k^s / (B_k^{s+p} \cap Z_k^s) = \frac{\text{Ker}(\partial_k^s)}{\text{Im}(\partial_{k+1}^{s+p}) \cup \text{Ker}(\partial_k^s)}.$$

159 The  $p$ -th persistent  $k$ -th Betti number  $\beta_k^{s,p}$  of  $K_s$  is the rank of  $H_k^{s,p}(K)$ . Note that the zero-persistent  
 160 homology groups of  $K_s$  are the same as the actual homology groups of  $K_s$ .

161 The results of computing the persistent homology of a filtered simplicial complex are normally given  
 162 in terms of persistent pairs consisting of birth times and death times. And, these are normally visualized  
 163 in various ways. The most common way of visualizing results of computing persistent homology is by  
 164 use of persistence barcodes. These are representations of the recorded birth times and death times of the  
 165 topological invariant properties or generators. Birth times and death times refer to the filtration values ( $\varepsilon$   
 166 values or filtration time) at which the generators appeared and vanished, respectively. The name persistent  
 167 barcode was first used in Zomorodian and Carlsson (2005) and discussed in more details in Ghrist (2008).  
 168 But, the standard algorithm for transforming persistent homology into barcodes is presented here as seen  
 169 in Otter et al. (2017), as their presentation is clear, concise and beginner friendly.

170 Given a filtered simplicial complex  $K$ , with filtration  $K_0 \subseteq K_1 \subseteq K_2 \subseteq \dots \subseteq K_r$ . The persistent pair  
 171 may be given by the pair  $(b, d)$ , where  $b \in \{0, 1, 2, \dots, r\}$  is the birth time,  $d \in (\{0, 1, 2, \dots, r\} \cup \infty)$  is the  
 172 death time,  $b \leq d$  and  $d - b$  is the length or lifespan of the homology. If  $d < \infty$  then the generator vanishes  
 173 at filtration time  $d$  and if  $d = \infty$  then the homology persists on all the succeeding filtration steps.

The persistent barcode for the filtered simplicial complex  $K$  can be created using the following steps.  
 First,  $K$  must be associated to boundary matrix whose entries represents faces of the simplexes. It is  
 assume that each of the simplexes of the nested sequence of complexes follow a total ordering such that a  
 face of a simplex precedes the simplex and a simplex in the  $i$ -th complex  $K_i$  precedes the simplices in  $K_j$   
 for  $j > i$ , which are not in  $K_i$ . Let  $n$  be the total number of simplices in the complex, and  $\sigma_1, \sigma_2, \dots, \sigma_n$   
 be the simplices. The square matrix  $B$ , of dimension  $n \times n$ , is constructed by assigning a value 1 in  $B(i, j)$  if  
 the simplex  $\sigma_i$  is a face of simplex  $\sigma_j$  of codimension 1 and a value 0 otherwise. That is, the boundary  
 matrix  $B$  is defined by

$$B(i, j) = \begin{cases} 1, & \text{if } \sigma_i \subset \sigma_j \text{ and } \dim(\sigma_j) - \dim(\sigma_i) = 1 \\ 0, & \text{otherwise.} \end{cases}$$

174 After constructing the boundary matrix  $B$ , it has to be reduced using the *standard algorithm*, sometimes  
 175 called column algorithm, for the computation of PH. The standard algorithm used for computing PH was  
 176 first introduced in Edelsbrunner et al. (2002). For each  $j \in \{1, 2, \dots, n\}$ , define  $\text{low}(j)$  to be the largest

177 value  $i$  such that  $B(i, j)$  is different from 0. If column  $j$  only contains 0 entries, then the value of  $low(j)$   
 178 is undefined. Boundary matrix  $B$  is reduced if the map  $low$  is injective on its domain of definition. In  
 179 the worst case, the complexity of the standard algorithm is cubic in the number of simplices (Otter et al.,  
 180 2017).

181 Finally, the reduced boundary matrix can now be encoded into a barcode. This is done by pairing the  
 182 simplices in the following manner:

- 183 • If  $low(j) = i$ , then the simplex  $\sigma_j$  is paired with  $\sigma_i$ , and the appearance of  $\sigma_i$  in the filtration causes  
 184 the birth of a feature that dies with the entrance of  $\sigma_j$ .
- 185 • If  $low(j)$  is undefined, then the appearance of the simplex  $\sigma_j$  in the filtration causes the birth of  
 186 a feature. If there exists  $k$  such that  $low(k) = j$ , then  $\sigma_j$  is paired with the simplex  $\sigma_k$ , whose  
 187 appearance in the filtration causes the death of the feature. If no such  $k$  exists, then  $\sigma_j$  is unpaired.

188 A pair  $(\sigma_i, \sigma_j)$  gives the half-open interval  $[dg(\sigma_i), dg(\sigma_j))$  in the barcode, where for a simplex  $\sigma \in K$   
 189 we define  $dg(\sigma)$  to be the smallest number  $l$  such that  $\sigma \in K_l$ . An unpaired simplex  $\sigma_k$  gives the infinite  
 190 interval  $[dg(\sigma_k), \infty)$ .

191 **Definition 12** Let  $K$  be a set of points in the Euclidean space  $\mathbb{R}^d$ . Fix  $m \in \mathbb{N}$  be the number of simplicial  
 192 complexes in the filtration of  $K$ . Suppose  $K^0 \subseteq K^1 \subseteq K^2 \subseteq \dots \subseteq K^m$  is a filtration of  $K$  with respect to  
 193 parameter values  $\varepsilon_i$ 's such that  $0 = \varepsilon_0 < \varepsilon_1 < \varepsilon_2 < \dots < \varepsilon_m$ . Let  $n \in \mathbb{N}$  be the number of simplices  $\sigma_j$ 's  
 194 in  $K^m$ . For each  $k \in \{0, 1, 2, \dots, d\}$ , there is a barcode  $B_k$  which is a collection of half-open intervals  
 195  $[\sigma_{j,1}, \sigma_{j,2})$ , which are pairs of birth time and death time of a generator in  $H_k(K)$ . Then, for each  $k$ ,  
 196  $k = 0, 1, 2, \dots, d$  the number of infinite intervals in  $B_k$  gives the  $k$ -th Betti number or the rank of  $H_k(K^m)$ .  
 197 And, the persistent homology of  $K$  based on the filtration  $K^0 \subseteq K^1 \subseteq K^2 \subseteq \dots \subseteq K^m$  is the homology of  
 198  $K^m$ .

199 **Example 5** Consider a set of points  $K = \{(0,0), (0,1), (.5, 1.5), (1,1), (1,0)\}$  and filtration values  $\varepsilon_1 =$   
 200  $0.2, \varepsilon_2 = 0.4, \varepsilon_3 = 0.6, \varepsilon_4 = 0.8, \text{ and } \varepsilon_5 = 1.0$ .

201 A filtration of  $K$  using the Čech complexes,  $\emptyset = K^0 \subseteq K^1 \subseteq K^2 \subseteq K^3 \subseteq K^4 \subseteq K^5$ , is given in Fig. 4.

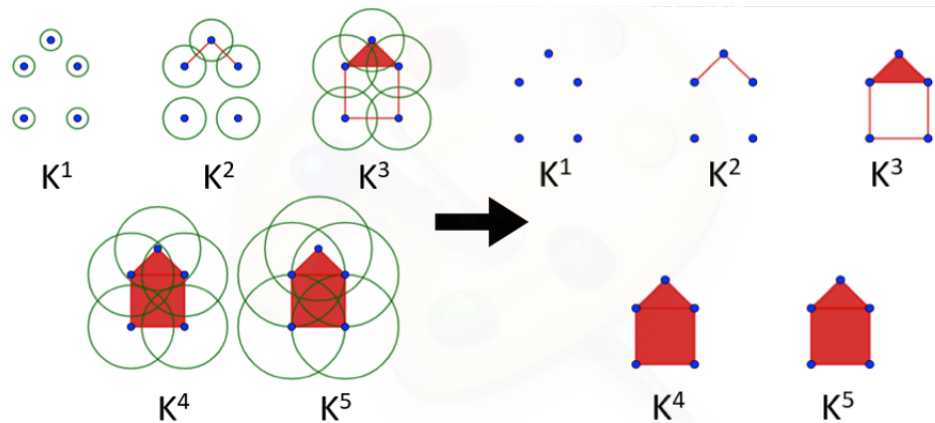


Figure 4. A Filtration of  $K$

202 After implementing the standard algorithm, its barcodes for degree 0 and degree 1 are shown in Fig. 5  
 203 and Fig. 6, respectively.

204 It can be concluded from the barcodes that the respective Betti numbers are  $\beta_0 = 1$  and  $\beta_1 = 0$ . Based  
 205 on the given filtration parameter values, the persistent homology of  $K$  can be described as a connected  
 206 space with no hole.

207 Persistence barcodes are then analyzed by studying properties of metric spaces whose elements are  
 208 persistence diagrams. A persistence diagram is another form of visualizing results of PH computations. It  
 209 gives similar information that a barcode provides. Also, distance functions were defined on a space of  
 210 persistence diagrams.

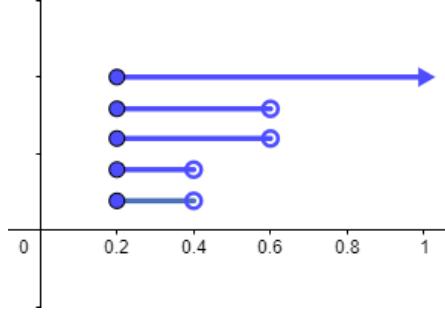


Figure 5. Barcode for degree 0

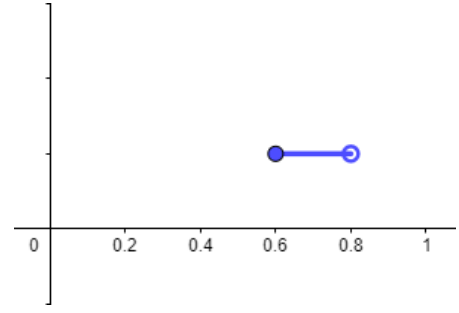


Figure 6. Barcode for degree 1

Recall that the  $n$ -th Betti number of a topological space  $X$  is denoted by  $\beta_n$ , which is equal to the rank of the  $n$ -th homology group  $H_n$ . Moreover, if  $K$  is a simplicial complex and  $\{K_r\}_{r \in J}$  for some indexing is the filtration of  $K$ , the  $p$ -th persistent  $k$ -th Betti number  $\beta_k^{s,p}$  of  $K_s$  is the rank of  $H_k^{s,p}(K)$ . From the persistent Betti numbers, there is a set of multiplicities  $\mu_n^{i,j} > i$  such that

$$p = j - i, \mu_n^{i,j} = \beta_n^{i,p} - \beta_n^{i-1,p} - \beta_n^{i,p+1} + \beta_n^{i-1,p+1}$$

211 The multiplicity  $\mu_n^{i,j}$  is the number of features in the  $n$ -th homology group that appears at filtration  $i$   
 212 and vanishes at filtration time  $j$ .

213 **Definition 13 (Persistence Diagram)** Let  $\{K_r\}$  be the filtration of a simplicial complex  $K$ . The  $n$ -th  
 214 persistence diagram of  $K$  with the filtration  $\{K_r\}$ , denoted by  $PD_n(\{K_r\})$  is a subset of  $\mathbb{R}^2$ , where  
 215  $\mathbb{R}^2 = (\mathbb{R} \cup \{\pm\infty\}) \times (\mathbb{R} \cup \{\pm\infty\})$ , with each point  $(i, j)$  has a multiplicity of  $\mu_n^{i,j}$  and all points in the  
 216 diagonal where  $i = j$  have infinite multiplicity.

217 Discussion of the robustness and stability of persistence diagram requires the notion of distance. Given  
 218 two persistence diagrams, say  $X$  and  $Y$ , the definition of distance between  $X$  and  $Y$  is given as follows.

**Definition 14** Let  $p \in [1, \infty]$ . The  $p$ -th Wasserstein distance between  $X$  and  $Y$  is defined as

$$W_p[d](X, Y) := \inf_{\phi: X \rightarrow Y} \left[ \sum_{x \in X} d[x, \phi(x)]^p \right]^{1/p}$$

for  $p \in [1, \infty)$  and as

$$W_\infty[d](X, Y) := \inf_{\phi: X \rightarrow Y} \sup_{x \in X} d[x, \phi(x)]$$

219 for  $p = \infty$ , where  $d$  is a metric on  $\mathbb{R}^2$  and  $\phi$  ranges over all bijections from  $X$  to  $Y$ .

220 Normally,  $d$  is taken to be  $L_q$  where  $q \in [1, \infty]$  and the most commonly used distance function is the  
 221 Bottleneck distance  $W_\infty[L_\infty]$ .

## 222 REFERENCES

- 223 Bubenik, P. (2015). Statistical topological data analysis using persistence landscapes. *Journal of Machine*  
 224 *Learning Research*, 16.  
 225 Carlsson, G. (2009). Topology and data. *Bull Am Math Soc*, 46.  
 226 Edelsbrunner, Letscher, and Zomorodian (2002). Topological persistence and simplification. *Discrete &*  
 227 *Computational Geometry*, 28(4):511–533.  
 228 Edelsbrunner, H. and Harer, J. (2008). Persistent homology — a survey.  
 229 Edelsbrunner, H. and Harer, J. (2010). *Computational topology: an introduction*. Am. Math. Soc.,  
 230 Providence.  
 231 Ghrist, R. (2008). Barcodes: the persistent topology of data. *Bull Am Math Soc*, 45.  
 232 Otter, N., Porter, M. A., Tillmann, U., Grindrod, P., and Harrington, H. A. (2017). A roadmap for the  
 233 computation of persistent homology. *EPJ Data Science*, 6(1):17.



- 234 Pun, C. S., Xia, K., and Lee, S. X. (2018). Persistent-homology-based machine learning and its applica-  
235 tions – a survey.
- 236 Zomorodian, A. and Carlsson, G. (2005). Computing persistent homology. *Discrete & Computational*  
237 *Geometry*, 33(2):249–274.