

PH-STAT

Moo K. Chung

University of Wisconsin-Madison mkchung@wisc.edu

Abstract. The PH-STAT toolbox performs various persistent homology based topological data analysis in MATLAB. The code and manual are distributed in <https://github.com/laplcebeltrami/ISBI2023TDA/tree/main/PH-STAT>.

1 Simplicial homology

A high dimensional object can be approximated by the point cloud data X consisting of p number of points. If we connect points of which distance satisfy a given criterion, the connected points start to recover the topology of the object. Hence, we can represent the underlying topology as a collection of the subsets of X that consists of nodes which are connected [2,5]. Given a point cloud data set X with a rule for connections, the topological space is a simplicial complex and its element is a simplex [12]. For point cloud data, the Delaunay triangulation is probably the most widely used method for connecting points. The Delaunay triangulation represents the collection of points in space as a graph whose face consists of triangles. Another way of connecting point cloud data is based on Rips complex often studied in persistent homology.

Homology is an algebraic formalism to associate a sequence of objects with a topological space [2]. In persistent homology, the algebraic formalism is usually built on top of objects that are hierarchically nested such as morse filtration, graph filtration and dendrograms. Formally homology usually refers to homology groups which are often built on top of a simplicial complex for point cloud and network data [7].

The k -simplex σ is the convex hull of $v + 1$ independent points v_0, \dots, v_k . A point is a 0-simplex, an edge is a 1-simplex, and a filled-in triangle is a 2-simplex. A *simplicial complex* is a finite collection of simplices such as points (0-simplex), lines (1-simplex), triangles (2-simplex) and higher dimensional counter parts [2]. A *k-skeleton* is a simplex complex of up to k simplices. Hence a graph is a 1-skeleton consisting of 0-simplices (nodes) and 1-simplices (edges). There are various simplicial complexes. The most often used simplicial complex in persistent homology is the Rips complex.

2 Rips complex

The Vietoris–Rips or Rips complex is the most often used simplicial complex in persistent homology. Let $X = \{x_0, \dots, x_p\}$ be the set of n points in \mathbb{R}^d . The distance matrix between points in X is given by $w = (w_{ij})$ where w_{ij} is the distance between points x_i and x_j . Then the Rips complex $R_\epsilon(X)$ is defined as follows [1,6]. The Rips complex is a collection of simplicial complexes parameterized by ϵ . The complex $R_\epsilon(X)$ captures the topology of the point set X at a scale of ϵ or less.

- The vertices of $R_\epsilon(X)$ are the points in X .

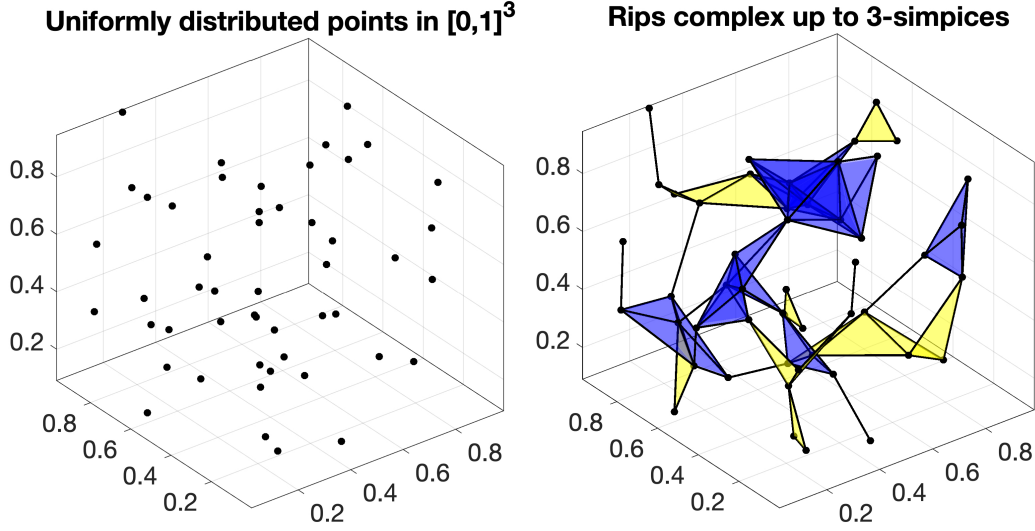


Fig. 1. Left: 50 randomly distributed points X in $[0, 1]^3$. Right: Rips complex $R_{0.3}(X)$ within radius 0.3 containing 106 1-simplices, 75 2-simplices (yellow) and 22 3-simplices (blue).

- If the distance w_{ij} is less than or equal to ϵ , then there is an edge connecting points x_i and x_j in $R_\epsilon(X)$.
- If the distance between any two points in $x_{i_0}, x_{i_1}, \dots, x_{i_k}$ is less than or equal to ϵ , then there is a k -simplex in $R_\epsilon(X)$ whose vertices are $x_{i_0}, x_{i_1}, \dots, x_{i_k}$.

In practice, the Rips complex is usually computed following the above definition iteratively adding simplices of increasing dimension. Given $p + 1$ number of points, there are potentially up to $\binom{p+1}{k}$ k -simplices making the data representation extremely inefficient as the radius ϵ increases. Thus, we restrict simplices of dimension up to k in practice. Such a simplicial complex is called the k -skeleton. It is implemented as `PH_rips.m`, which inputs the matrix X of size $p \times d$, dimension k and radius e . Then outputs the structured array S containing the collection of nodes, edges, faces up to k -simplices. For instance, the Rips complex up to 3-simplices in Figure 1 is created using

```
p=50; d=3;
X = rand(p, d);
S= PH_rips(X, 3, 0.3)
PH_rips_display(X,S);
```

$S =$

4×1 cell array

```
{ 50×1 double}
{106×2 double}
```

```
{ 75×3 double}
{ 22×4 double}
```

The Rips complex is then displayed using `PH_rips_display.m` which inputs node coordinates \mathbf{X} and simplicial complex \mathbf{S} .

3 Boundary matrix

Given a simplicial complex K , the boundary matrices B_k represent the boundary operators between the simplices of dimension k and $k - 1$. Let C_k be the collection of k -simplices. Define the k -th boundary map

$$\partial_k : C_k \rightarrow C_{k-1}$$

as a linear map that sends each k -simplex σ to a linear combination of its $k - 1$ faces

$$\partial_k \sigma = \sum_{\tau \in F_k(\sigma)} (-1)^{\text{sgn}(\tau, \sigma)} \tau,$$

where $F_k(\sigma)$ is the set of $k - 1$ faces of σ , and $\text{sgn}(\tau, \sigma)$ is the sign of the permutation that sends the vertices of τ to the vertices of σ . This expression says that the boundary of a k -simplex σ is the sum of all its $(k - 1)$ -dimensional faces, with appropriate signs determined by the orientation of the faces. The signs alternate between positive and negative depending on the relative orientation of the faces, as determined by the permutation that maps the vertices of one face to the vertices of the other face. The k -th boundary map removes the filled-in interior of k -simplices.

Consider a filled-in triangle $\sigma = [v_1, v_2, v_3] \in C_2$ with three vertices v_1, v_2, v_3 in Figure 2. The boundary map ∂_k applied to σ resulted in the collection of three edges that forms the boundary of σ :

$$\partial_2 \sigma = [v_1, v_2] + [v_2, v_3] + [v_3, v_1] \in C_1. \quad (1)$$

If we give the direction or orientation to edges such that

$$[v_3, v_1] = -[v_1, v_3],$$

and use edge notation $e_{ij} = [v_i, v_j]$, we can write (1) as

$$\partial_2 \sigma = e_{12} + e_{23} + e_{31} = e_{12} + e_{23} - e_{13}.$$

The boundary map can be represented as a boundary matrix ∂_k with respect to a basis of the vector spaces C_k and C_{k-1} , where the rows of ∂_k correspond to the basis elements of C_k and the columns correspond to the basis elements of C_{k-1} . The (i, j) entry of ∂_k is given by the coefficient of the j th basis element in the linear combination of the $k - 1$ faces of the i th basis element in C_k . The boundary matrix is the higher dimensional version of the incidence matrix in graphs [9,8,11] showing how $(k - 1)$ -dimensional simplices are forming k -dimensional simplex. The (i, j) entry of ∂_k is one if τ is a face of σ otherwise zero. The entry can be -1 depending on the orientation of τ . For the simplicial complex in Figure 2,

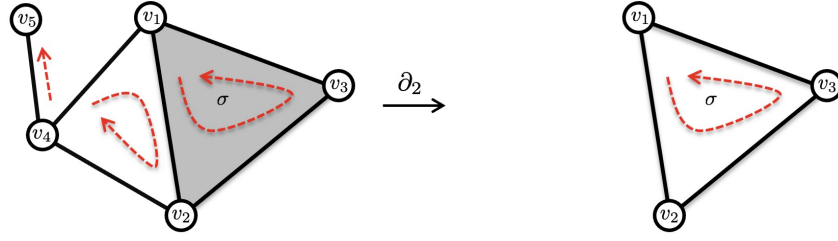


Fig. 2. A simplicial complex with 5 vertices and 2-simplex $\sigma = [v_1, v_2, v_3]$ with a filled-in face (colored gray). After boundary operation ∂_2 , we are only left with 1-simplices $[v_1, v_2] + [v_2, v_3] + [v_3, v_1]$, which is the boundary of the filled in triangle. The complex has a single connected component ($\beta_0 = 1$) and a single 1-cycle. The dotted red arrows are the orientation of simplices.

the boundary matrices are given by

$$\partial_2 = \begin{matrix} & \sigma \\ e_{12} & \begin{pmatrix} 1 \\ 1 \\ 1 \\ 0 \\ 0 \\ 0 \end{pmatrix} \\ e_{23} & \\ e_{31} & \\ e_{24} & \\ e_{41} & \\ e_{45} & \end{matrix}$$

$$\partial_1 = \begin{matrix} & e_{12} & e_{23} & e_{31} & e_{24} & e_{41} & e_{45} \\ v_1 & \begin{pmatrix} -1 & 0 & 1 & 0 & 1 & 0 \end{pmatrix} \\ v_2 & \begin{pmatrix} 1 & -1 & 0 & -1 & 0 & 0 \end{pmatrix} \\ v_3 & \begin{pmatrix} 0 & 1 & -1 & 0 & 0 & 0 \end{pmatrix} \\ v_4 & \begin{pmatrix} 0 & 0 & 0 & 1 & -1 & -1 \end{pmatrix} \\ v_5 & \begin{pmatrix} 0 & 0 & 0 & 0 & 0 & 1 \end{pmatrix} \end{matrix}$$

$$\partial_0 = \begin{matrix} & v_1 & v_2 & v_3 & v_4 & v_5 \\ 0 & \begin{pmatrix} 0 & 0 & 0 & 0 & 0 \end{pmatrix} \end{matrix}.$$

In example in Figure 3-left, `PH_rips(X,3,0.5)` gives

```
>> S{1}
```

```
1
2
3
4
5
```

```
>> S{2}
```

```
1    3
2    5
3    4
```

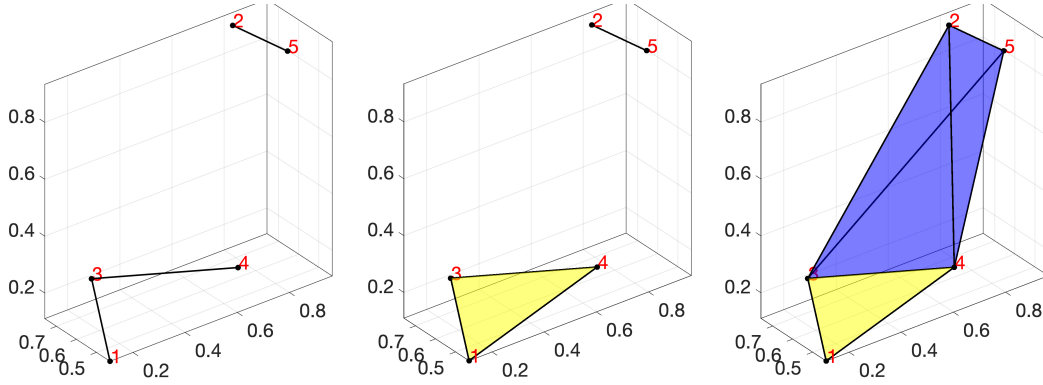


Fig. 3. Examples of boundary matrix computation. From the left to right, the radius is changed to 0.5, 0.6 and 1.0.

PH_boundary.m only use node set $S\{1\}$ and edge set $S\{2\}$ in building boundary matrix $B1$ saving computer memory.

```
>> B{1}
```

```
-1    0    0
 0   -1    0
 1    0   -1
 0    0    1
 0    1    0
```

The columns of boundary matrix $B\{1\}$ is indexed with the edge set in $S\{2\}$ such that the first column corresponds to edge $[1,3]$. Any other potential edges $[2,3]$ that are not connected is simply ignored to save computer memory.

When we increase the filtration value and compute $PH_rips(X,3,0.6)$, a triangle is formed (yellow colored) and $S\{3\}$ is created (Figure 3-middle).

```
>> S{2}
```

```
1    3
1    4
2    5
3    4
```

```
>> S{3}
```

```
1    3    4
```

Correspondingly, the boundary matrices change to

```
>> B{1}
```

-1	-1	0	0
0	0	-1	0
1	0	0	-1
0	1	0	1
0	0	1	0

>> B{2}

1
-1
0
1

From the edge set $S\{2\}$ that forms the row index for boundary matrix $B\{2\}$, we have $[1, 3] - [1, 4] + [3, 4]$ that forms the triangle $[1, 3, 4]$.

When we increase the filtration value further and compute $\text{PH_rips}(X, 3, 1)$, a tetrahedron is formed (blue colored) and $S\{4\}$ is created (Figure 3-right).

>> S{3}

1	3	4
2	3	4
2	3	5
2	4	5
3	4	5

>> S{4}

2	3	4	5
---	---	---	---

Correspondingly, the boundary matrix $B\{3\}$ is created

>> B{3}

0
-1
1
-1
1

The easiest way to check the computation is correct is looking at the sign of triangles in $-[2, 3, 4] + [2, 3, 5] - [2, 4, 5] + [3, 4, 5]$. Using the right hand thumb rule, which puts the orientation of triangle $[3, 4, 5]$ toward the center of the tetrahedron, the orientation of all the triangles are toward the center of the tetrahedron. Thus, the signs are correctly assigned. Since computer algorithms are built inductively, the method should work correctly in higher dimensional simplices.

4 Homology group

The boundary map satisfy the property that the composition of any two consecutive boundary maps is the zero map, i.e.,

$$\partial_{k-1} \circ \partial_k = 0.$$

This reflect the fact that the boundary of a boundary is always empty. This property implies that the image of ∂_k is contained in the kernel of ∂_{k-1} . The kernel of boundary map is defined as

$$\ker \partial_k = \{\sigma \in C_k | \partial_k \sigma = 0\}.$$

The elements of the kernel of ∂_k are called k -cycles, since they form closed loops or cycles in the simplicial complex. The kernel of the boundary matrix B_k is spanned by eigenvectors v corresponding to zero eigenvalues of B_k .

We can apply the boundary operation ∂_1 further to $\partial_2 \sigma$ and obtain

$$\begin{aligned} \partial_1 \partial_2 \sigma &= \partial_1 e_{12} + \partial_1 e_{23} + \partial_1 e_{31} \\ &= v_2 - v_1 + v_3 - v_2 + v_1 - v_3 = 0. \end{aligned}$$

The boundary operation twice will results in an empty set. Such algebraic representation for boundary operation has been very useful for effectively quantifying persistent homology.

The image of boundary map is defined as

$$\text{im} \partial_{k+1} = \{\partial_{k+1} \sigma | \sigma \in C_{k+1}\}.$$

The elements of the image of ∂_{k+1} are called k -boundaries, and they represent k -dimensional features that can be filled in by $(k+1)$ -dimensional simplices. The image of the boundary matrix B_{k+1} is the subspace spanned by its columns. The column space can be found by the Gaussian elimination or singular value decomposition.

We can define the k th homology group $H_k(K)$ as the quotient space of the kernel of ∂_k modulo the image of ∂_{k+1} :

$$H_k(K) = \ker(\partial_k) / \text{im}(\partial_{k+1}).$$

$H_k(K)$ is a vector space that captures the k th topological feature or hole in K .

[Given boundary matrix $S\{d+1\}$ and $S\{d\}$, compute and represent the quotient space in Matlab]

Intuitively, it measures the presence of k -dimensional loops in the simplicial complex. The rank of $H_k(K)$ is the k th Betti number of K , which is an algebraic invariant that captures the topological features of the complex K .

Although we put direction in the boundary matrix ∂_1 by adding sign, the Betti number β_1 computation will be invariant. With boundary operations, we can build a vector space C_k using the set of k -simplices as a basis. The vector spaces $C_k, C_{k-1}, C_{k-2}, \dots$ are then sequentially nested by boundary operator ∂_k [2]:

$$\dots \xrightarrow{\partial_{k+1}} C_k \xrightarrow{\partial_k} C_{k-1} \xrightarrow{\partial_{k-1}} C_{k-2} \xrightarrow{\partial_{k-2}} \dots \quad (2)$$

Such nested structure is called the *chain complex*. Let B_k be a collection of boundaries obtained as the image of ∂_k , i.e.,

$$B_k = \{\partial_k \sigma : \sigma \in C_k\}.$$

Let Z_k be a collection of cycles obtained as the kernel of ∂_k , i.e.,

$$Z_k = \{\sigma \in C_k : \partial_k \sigma = 0\}.$$

For instance, the 1-cycle formed by edges e_{12}, e_{23}, e_{31} in Figure 2 is the boundary of the filled-in gray colored triangle σ . The boundaries B_k form subgroups of the cycles Z_k , i.e., $B_k \subset Z_k$. We can partition Z_k into cycles that differ from each other by boundaries through the quotient space

$$H_k = Z_k / B_k,$$

which is called the k -th homology group. The elements of the k -th homology group are often referred to as k -dimensional cycles or k -cycles. The k -th Betti number β_k is then the number of k -dimensional cycles, which is given by the rank of H_k , i.e.,

$$\beta_k = \text{rank}(H_k) = \text{rank}(Z_k) - \text{rank}(B_k). \quad (3)$$

The 0-th Betti number is the number of connected components while the 1-st Betti number is the number of cycles.

[Given boundary matrix $S\{d+1\}$ and $S\{d\}$, compute the Betti numbers]

The Betti numbers β_k are usually algebraically computed by reducing the boundary matrix ∂_k to the Smith normal form, which has a block diagonal matrix as a submatrix in the upper left, via Gaussian elimination [2]. For instance, the boundary matrices ∂_i in Figure 2 is transformed to the Smith normal form $\mathcal{S}(\partial_i)$ after Gaussian elimination as

$$\mathcal{S}(\partial_1) = \begin{pmatrix} 1 & 0 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 0 & 0 \end{pmatrix}, \quad \mathcal{S}(\partial_2) = \begin{pmatrix} 1 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \end{pmatrix}.$$

In the Smith normal form $\mathcal{S}(\partial_k)$, the number of columns containing only zeros is $\text{rank}(Z_k)$, the number of k -cycles while the number of rows containing one is $\text{rank}(B_{k-1})$, the number of $(k-1)$ -cycles that are boundaries. From (3), the Betti number computation involves the rank computation of two boundary matrices. In Figure 2 example, there are $\text{rank}(Z_1) = 2$ zero columns and $\text{rank}(B_0) = 4$ non-zero rows. $\text{rank}(Z_0) = 5$ is trivially the number of nodes in the simplicial complex while there are $\text{rank}(B_1) = 1$ for $\mathcal{S}(\partial_2)$. Thus, we have

$$\begin{aligned} \beta_0 &= \text{rank}(Z_0) - \text{rank}(B_0) = 5 - 4 = 1, \\ \beta_1 &= \text{rank}(Z_1) - \text{rank}(B_1) = 2 - 1 = 1. \end{aligned}$$

The Betti numbers can be also computed using the Hodge Laplacian without Gaussian elimination. The standard graph Laplacian is defined as

$$\Delta_0 = \partial_1 \partial_1^\top,$$

which is also called the 0-th Hodge Laplacian [9]. In general, the k -th Hodge Laplacian is defined as

$$\Delta_k = \partial_{k+1} \partial_{k+1}^\top + \partial_k \partial_k^\top.$$

The boundary operation ∂_k only depends on k -simplices. Thus, Δ_k is uniquely determined by $(k+1)$ - and k -simplices. The k -th Laplacian is sparse a $n_k \times n_k$ positive semi-definite

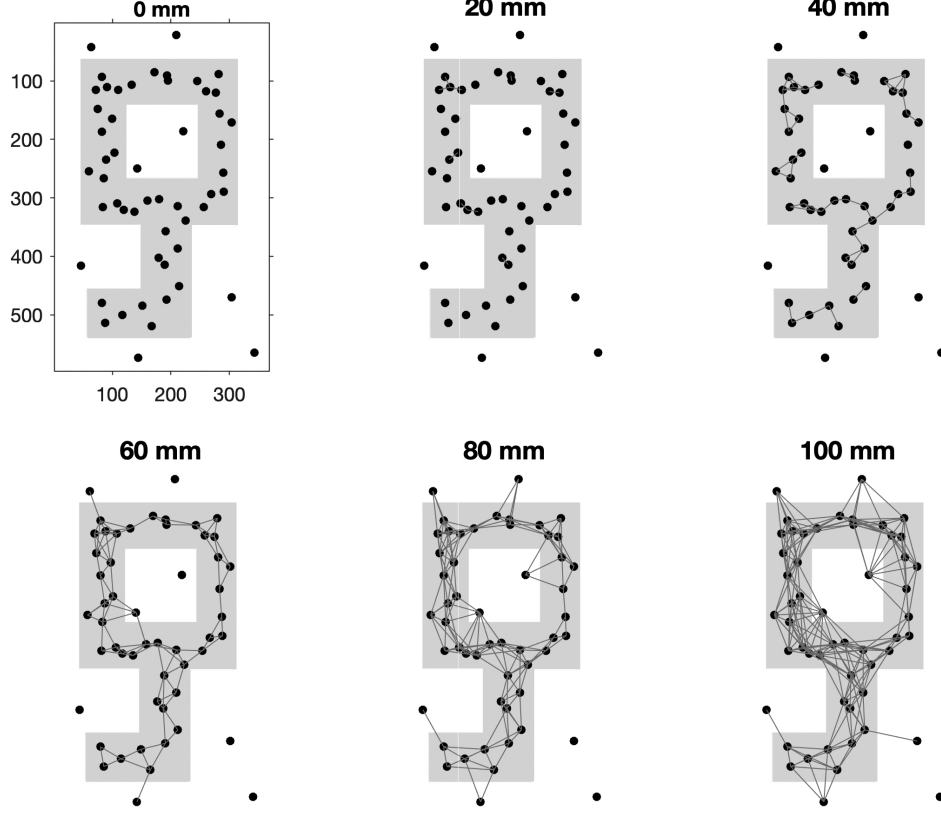


Fig. 4. Rips filtration on 1-skeleton (Rips complex consisting of only nodes and edges) of the point cloud data that was sampled along the underlying key shaped data. If two points are within the given radius, we connect them with an edge.

symmetric matrix, where n_k is the number of k -simplices in the network [3]. Then the k -th Betti number β_k is the dimension of $\ker \Delta_k$, which is given by computing the rank of Δ_k . The 0th Betti number, the number of connected component, is computed from Δ_0 while the 1st Betti number, the number of cycles, is computed from Δ_1 .

The k -th hodge Laplacian depends only on n_k number of k -simplices in the data. After lengthy algebraic derivation, we can show that

$$\Delta_k = D - A_u + (k + 1)I_{n_k} + A_l,$$

where A_u and A_l are the upper and lower adjacency matrices between the k -simplices. $D = \text{diag}(\deg(\sigma_1), \dots, \deg(\sigma_{n_k}))$ is the diagonal matrix consisting of the sum of node degrees of simplices σ_j [10].

4.1 Rips filtrations: filtrations on point cloud data

The Rips complex has been the main building block for persistent homology and defined on top of the point cloud data [4]. The *Rips complex* is a simplicial complex constructed by

connecting two data points if they are within specific distance ϵ . Figure 4 shows an example of the Rips complex that approximates the gray object with a point cloud. Given a point cloud data, the Rips complex R_ϵ is a simplicial complex whose k -simplices correspond to unordered $(k + 1)$ -tuples of points which are pairwise within distance ϵ [4]. While a graph has at most 1-simplices, the Rips complex has at most k -simplices. The Rips complex has the property that

$$\mathcal{R}_{\epsilon_0} \subset \mathcal{R}_{\epsilon_1} \subset \mathcal{R}_{\epsilon_2} \subset \dots$$

for $0 = \epsilon_0 \leq \epsilon_1 \leq \epsilon_2 \leq \dots$. When $\epsilon = 0$, the Rips complex is simply the node set V . By increasing the filtration value ϵ , we are connecting more nodes so the size of the edge set increases. Such the nested sequence of the Rips complexes is called a *Rips filtration*, the main object of interest in the persistent homology [1]. The increasing ϵ values are called the filtration values.

One major problem of the Rips complex is that as the number of vertices p increase, the resulting simplicial complex becomes very dense. Further, as the filtration values increases, there exists an edge between ever pair of vertices and filled triangle between every triple of vertices. At higher filtration values, Rips filtration becomes very ineffective representation of data.

References

1. Edelsbrunner, H., Harer, J.: Persistent homology - a survey. Contemporary Mathematics **453**, 257–282 (2008)
2. Edelsbrunner, H., Harer, J.: Computational topology: An introduction. American Mathematical Society (2010)
3. Friedman, J.: Computing betti numbers via combinatorial laplacians. Algorithmica **21**(4), 331–346 (1998)
4. Ghrist, R.: Barcodes: The persistent topology of data. Bulletin of the American Mathematical Society **45**, 61–75 (2008)
5. Hart, J.: Computational topology for shape modeling. In: Proceedings of the International Conference on Shape Modeling and Applications. pp. 36–43 (1999)
6. Hatcher, A.: Algebraic topology. Cambridge University Press (2002)
7. Lee, H., Chung, M.K., K.H., Lee, D.: Hole detection in metabolic connectivity of Alzheimer’s disease using k-Laplacian. MICCAI, Lecture Notes in Computer Science **8675**, 297–304 (2014)
8. Lee, H., Chung, M., Choi, H., K., H., Ha, S., Kim, Y., Lee, D.: Harmonic holes as the submodules of brain network and network dissimilarity. International Workshop on Computational Topology in Image Context, Lecture Notes in Computer Science pp. 110–122 (2019)
9. Lee, H., Chung, M., Kang, H., Choi, H., Kim, Y., Lee, D.: Abnormal hole detection in brain connectivity by kernel density of persistence diagram and Hodge Laplacian. In: IEEE International Symposium on Biomedical Imaging (ISBI). pp. 20–23 (2018)
10. Muhammad, A., Egerstedt, M.: Control using higher order laplacians in network topologies. In: Proc. of 17th International Symposium on Mathematical Theory of Networks and Systems. pp. 1024–1038 (2006)
11. Schaub, M., Benson, A., Horn, P., Lippner, G., Jadbabaie, A.: Random walks on simplicial complexes and the normalized hodge laplacian. arXiv preprint arXiv:1807.05044 (2018)
12. Zomorodian, A.: Topology for computing. Cambridge University Press, Cambridge (2009)