

Topological Data Analysis in Python

organized by: Michael Bleher, Maximilian Schmahl and Daniel Spitz

Heidelberg University

26th - 28th of October 2020

Contents

Homology of Spaces

Homology of Data

Ripser

Applications

Contents

Homology of Spaces

Homology of Data

Ripser

Applications

Algebraic Topology: Distinguishing topological spaces via algebraic invariants

Algebraic Topology: Distinguishing topological spaces via algebraic invariants Homology: For each dimension $d \ge 0$ assigns to a topological space a vector space

 $X \mapsto H_d(X)$

and to a continuous map a linear map

 $(f: X \to Y) \mapsto (H_d(f): H_d(X) \to H_d(Y))$

Algebraic Topology: Distinguishing topological spaces via algebraic invariants Homology: For each dimension $d \ge 0$ assigns to a topological space a vector space

 $X \mapsto H_d(X)$

and to a continuous map a linear map

$$(f: X \to Y) \mapsto (H_d(f): H_d(X) \to H_d(Y))$$

Definition $\beta_d(X) = \dim H_d(X)$ is called the *d*-th Betti number of X.

What homology can distinguish I



What homology can distinguish II



What homology can distinguish III



What homology cannot distinguish I



What homology cannot distinguish II



$\beta_d(X)$ is the number of *d*-dimensional holes in *X*.

A geometric *n*-simplex is the convex hull of n + 1 affinely independent points in \mathbb{R}^m .

. r ~

A geometric *n*-simplex is the convex hull of n + 1 affinely independent points in \mathbb{R}^m . If conv X is a simplex and $Y \subseteq X$, then conv Y is called a *face* of conv X.



A geometric *n*-simplex is the convex hull of n + 1 affinely independent points in \mathbb{R}^m . If conv X is a simplex and $Y \subseteq X$, then conv Y is called a *face* of conv X.

A geometric simplicial complex is a finite set K of simplices such that

- ▶ if $\sigma \in K$ and τ is a face of K, then $\tau \in K$, and
- if $\sigma, \tau \in K$ and $\sigma \cap \tau \neq \emptyset$, then $\sigma \cap \tau$ is a face of σ and of τ .



Boundary matrix

Definition Let $K = \{\sigma_1, \ldots, \sigma_m\}$ be a simplicial complex. We define its boundary matrix $D = (d_{i,j})$ as $d_{i,j} = \begin{cases} 1 & \text{if } \sigma_i \text{ is a boundary of } \sigma_j \text{ with } \dim \sigma_i = \dim \sigma_j - 1, \\ 0 & \text{otherwise.} \end{cases}$

ᢧ

011

000

Definition

Let $K = \{\sigma_1, \ldots, \sigma_m\}$ be a simplicial complex. We define its *boundary matrix* $D = (d_{i,j})$ as

$$d_{i,j} = \begin{cases} 1 & \text{if } \sigma_i \text{ is a boundary of } \sigma_j \text{ with } \dim \sigma_i = \dim \sigma_j - 1, \\ 0 & \text{otherwise.} \end{cases}$$

Example

ple
$$\int_{-1}^{1} \int_{-1}^{1} \int_{-1$$

Python Course on Topological Methods in Data Analysis - Day 2

Definition Let $v = (v_1, ..., v_m)^T$ be a column vector. We define pivot $v = \max\{i \in \{1, ..., m\} \mid v_i \neq 0\}.$

1

Definition Let $v = (v_1, ..., v_m)^T$ be a column vector. We define

pivot
$$v = \max\{i \in \{1, ..., m\} \mid v_i \neq 0\}.$$

If M is a matrix, we write m_j for its j-th column and cols M for the set of all its non-zero columns.

Definition Let $v = (v_1, ..., v_m)^T$ be a column vector. We define

pivot
$$v = \max\{i \in \{1, ..., m\} \mid v_i \neq 0\}.$$

If *M* is a matrix, we write m_j for its *j*-th column and cols *M* for the set of all its non-zero columns. We write pixets $M = \{p_i | m_i \in cols M\}$

We write pivots $M = \{ \text{pivot } m_j \mid m_j \in \text{cols } M \}.$

Definition Let $v = (v_1, ..., v_m)^T$ be a column vector. We define

pivot
$$v = \max\{i \in \{1, ..., m\} \mid v_i \neq 0\}.$$

If M is a matrix, we write m_j for its j-th column and cols M for the set of all its non-zero columns.

We write pivots $M = \{ \text{pivot } m_j \mid m_j \in \text{cols } M \}$. We say that M is *reduced* if pivot $m_j = p_{k}$ implies $m_j = m_k$ for all $m_j, m_k \in \text{cols } M$.

Pivotm

Definition Let $v = (v_1, ..., v_m)^T$ be a column vector. We define

pivot
$$v = \max\{i \in \{1, ..., m\} \mid v_i \neq 0\}.$$

If M is a matrix, we write m_j for its j-th column and cols M for the set of all its non-zero columns.

We write pivots $M = \{ \text{pivot } m_j \mid m_j \in \text{cols } M \}$. We say that M is *reduced* if pivot $m_j = m_k$ implies $m_j = m_k$ for all $m_j, m_k \in \text{cols } M$.

Example

Reduced:

Not reduced:

$$\begin{pmatrix}
0 & 7 & 1 \\
1 & 0 & 0 \\
0 & 0 & 1
\end{pmatrix}$$

$$\begin{pmatrix}
0 & 0 & 1 \\
1 & 0 & 7 \\
2 & 0 & 1
\end{pmatrix}$$

Maximilian Schmahl

Let *M* be a matrix. We say that (R, V) is a *reduction* of *M* if *R* is reduced, *V* is upper-triangular and invertible and we have R = MV.

```
Let M be a matrix. We say that (R, V) is a reduction of M if R is reduced, V is upper-triangular and invertible and we have R = MV.
```

Input: MOutput: (R, V) reduction of M $R \leftarrow M$ $V \leftarrow I$ while $\exists i < j$ with pivot $R_i = \text{pivot } R_j$ do $R_j \leftarrow R_j + R_i$ $V_j \leftarrow V_j + V_i$ end while return (R, V)

$$\begin{pmatrix} 1 & 0 & 1 \\ 0 & 1 & 1 \\ 0 & 1 & 1 \end{pmatrix} \longrightarrow \begin{pmatrix} 1 & 0 & 1 \\ 0 & 1 & 0 \\ 0 & 1 & 0 \end{pmatrix} \longrightarrow \begin{pmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 1 & 0 \end{pmatrix}$$

Theorem

Let $K = \{\sigma_1, \ldots, \sigma_m\}$ be a simplicial complex with boundary matrix D. Let (R, V) be a reduction of D. Then

$$\beta_d(K) = \#\{i \mid R_i = 0, i \notin \text{pivots } R, \dim \sigma_i = d\}.$$

Theorem

Let $K = \{\sigma_1, \ldots, \sigma_m\}$ be a simplicial complex with boundary matrix D. Let (R, V) be a reduction of D. Then

 $\beta_d(K) = \#\{i \mid R_i = 0, i \notin \text{pivots } R, \dim \sigma_i = d\}.$

Computing homology II



Homology of Spaces

Homology of Data

Ripser

Applications











Maximilian Schmahl

Python Course on Topological Methods in Data Analysis - Day 2

Theorem (Niyogi, Smale, Weinberger) Let $M \subseteq \mathbb{R}^k$ be a manifold and $P \subseteq M$. There exists c(M) > 0 such that for every $c(M) > \delta > 0$ with $M \subseteq \bigcup_{x \in P} B_{\delta}(x)$ we have $H_*(M) \cong H_*\left(\bigcup_{x \in P} B_{\delta}(x)\right)$

Theorem (Niyogi, Smale, Weinberger) Let $M \subseteq \mathbb{R}^k$ be a manifold and $P \subseteq M$. There exists c(M) > 0 such that for every $c(M) > \delta > 0$ with $M \subseteq \bigcup_{x \in P} B_{\delta}(x)$ we have

$$H_*(M) \cong H_*\left(\bigcup_{x\in P} B_\delta(x)\right)$$

How do we choose δ ?
Theorem (Niyogi, Smale, Weinberger) Let $M \subseteq \mathbb{R}^k$ be a manifold and $P \subseteq M$. There exists c(M) > 0 such that for every $c(M) > \delta > 0$ with $M \subseteq \bigcup_{x \in P} B_{\delta}(x)$ we have

$$H_*(M) \cong H_*\left(\bigcup_{x\in P} B_\delta(x)\right)$$

How do we choose δ ? We don't!









Maximilian Schmahl

Definition Let $P \subseteq \mathbb{R}^k$ be a finite set. For t > 0 define the *Čech complex* $\check{\mathsf{C}}\mathsf{ech}_t(P) = \{ Q \subseteq P \mid \bigcap B_t(x) \neq \emptyset \}.$ $x \in Q$ Note: $\check{C}ech_t(P) \subseteq \check{C}ech_u(P)$ whenever $t \leq u$.

Maximilian Schmahl

Definition Let $P \subseteq \mathbb{R}^k$ be a finite set. For t > 0 define the *Čech complex*

$$\check{\operatorname{Cech}}_t(P) = \{ Q \subseteq P \mid \bigcap_{x \in Q} B_t(x) \neq \emptyset \}.$$

Note: $\check{C}ech_t(P) \subseteq \check{C}ech_u(P)$ whenever $t \leq u$.

Theorem (Nerve Theorem)

Let $P \subseteq \mathbb{R}^k$ be a finite set and $t \ge 0$. If $\bigcap_{x \in Q} B_t(x)$ can be deformed to a point for all $Q \subseteq P$, then

$$H_*\left(\bigcup_{x\in P}B_t(x)\right)\cong H_*(\check{\operatorname{Cech}}_t(P))$$

Rips complexes

Definition Let $P \subseteq \mathbb{R}^k$ be a finite set. For $t \ge 0$ we define the *Rips complex*

$$\operatorname{Rips}_t(P) = \{ Q \subseteq P \mid \sup_{x,y \in Q} d(x,y) \leq t \}.$$

Note: $\operatorname{Rips}_t(P) \subseteq \operatorname{Rips}_u(P)$ whenever $t \leq u$.



Rips complexes

Definition Let $P \subseteq \mathbb{R}^k$ be a finite set. For $t \ge 0$ we define the *Rips complex*

$$\operatorname{Rips}_t(P) = \{ Q \subseteq P \mid \sup_{x,y \in Q} d(x,y) \leq t \}.$$

Note: $\operatorname{Rips}_t(P) \subseteq \operatorname{Rips}_u(P)$ whenever $t \leq u$.

Theorem

There exists $\theta > 0$ such that for all finite sets $P \subseteq \mathbb{R}^k$ and t > 0 we have

$$\mathsf{Rips}_t \subseteq \check{\mathsf{Cech}}_{\theta t}(P)$$

and

$$\check{\mathrm{Cech}}_t \subseteq \mathrm{Rips}_{\theta t}(P)$$





Maximilian Schmahl



Maximilian Schmahl



Maximilian Schmahl



Maximilian Schmahl



Maximilian Schmahl



Maximilian Schmahl

Barcodes

Definition

A family of finite-dimensional vector spaces V_i , i = 1, ..., n with linear maps $V_i \rightarrow V_{i+1}$ for each i is called a *persistence module*.

Barcodes

$$K_1 \subseteq K_2 \subseteq K_1 \subseteq \dots \subseteq K_n$$

Definition $\mathcal{H}(\mathcal{K}_{n}) \longrightarrow \mathcal{H}(\mathcal{K}_{n}) \longrightarrow \mathcal{H}(\mathcal{K}_{n})$ A family of finite-dimensional vector spaces V_{i} , i = 1, ..., n with linear maps $V_{i} \rightarrow V_{i+1}$ for each i is called a *persistence module*.

Theorem Gabriel
$$(V_{1})_{tell}$$
 $V_{1} \rightarrow V_{s}$ set
 $V_{1} \longrightarrow V_{2} \longrightarrow \dots \longrightarrow V_{n-1} \longrightarrow V_{n}$

be a persistence module. Then there exists a unique family of intervals $(I_k)_{k \in K}$ with $I_k \subseteq \{1, \ldots, n\}$ such that

$$\dim V_i = \#\{k \in K \mid i \in I_k\}$$

and

$$\operatorname{rank}(V_i \to V_j) = \#\{k \in K \mid i, j \in I_k\}$$

Maximilian Schmahl

Definition

Let K be a simplicial complex. A *filtration* of K is a family of simplicial complexes K_1, \ldots, K_n such that $K_n = K$ and $K_i \subseteq K_{i+1}$ for all n.

Definition

Let K be a simplicial complex. A *filtration* of K is a family of simplicial complexes K_1, \ldots, K_n such that $K_n = K$ and $K_i \subseteq K_{i+1}$ for all n. We call

$$H_*(K_1) \longrightarrow H_*(K_2) \longrightarrow \ldots \longrightarrow H_*(K_{n-1}) \longrightarrow H_*(K_n)$$

the *persistent homology* of the filtration.

Computing persistent homology

Theorem (Barannikov, Carlsson & Zomorodian, Edelsbrunner et al.) σ_{3} Let $K = \{\sigma_{1}, \ldots, \sigma_{n}\}$ be a simplicial complex such that $K_{i} = \{\sigma_{1}, \ldots, \sigma_{i}\}$ is a subcomplex.

Computing persistent homology

Theorem (Barannikov, Carlsson & Zomorodian, Edelsbrunner et al.) Let $K = \{\sigma_1, \ldots, \sigma_n\}$ be a simplicial complex such that $K_i = \{\sigma_1, \ldots, \sigma_i\}$ is a subcomplex.

Let D be the corresponding boundary matrix and (R, V) a reduction of D.

Computing persistent homology

Theorem (Barannikov, Carlsson & Zomorodian, Edelsbrunner et al.) Let $K = \{\sigma_1, \ldots, \sigma_n\}$ be a simplicial complex such that $K_i = \{\sigma_1, \ldots, \sigma_i\}$ is a subcomplex.

Let D be the corresponding boundary matrix and (R, V) a reduction of D. Then

$$\{[i,\infty) \mid R_i = 0, i \notin \text{pivots } R\} \cup \{[i,j) \mid i = \text{pivot } R_j\}$$

is the barcode of

$$H_*(K_1) \longrightarrow H_*(K_2) \longrightarrow \ldots \longrightarrow H_*(K_{n-1}) \longrightarrow H_*(K_n)$$

Persistent homology pipeline



Definition Let $P, Q \subseteq \mathbb{R}^k$. We define their *Hausdorff distance* as

$$d_H(P,Q) = \max \left\{ \sup_{p \in P} \inf_{q \in Q} d(p,q), \sup_{q \in Q} \inf_{p \in P} d(p,q) \right\}$$

Definition Let $P, Q \subseteq \mathbb{R}^k$. We define their *Hausdorff distance* as

$$d_H(P,Q) = \max \left\{ \sup_{p \in P} \inf_{q \in Q} d(p,q), \sup_{q \in Q} \inf_{p \in P} d(p,q) \right\}$$

Example



Definition Let $B = (I_{\alpha})_{\alpha \in A}$ and $C = (J_{\beta})_{\beta \in B}$ be barcodes. A δ -matching between them consists of subsets $A' \subseteq A$, $B' \subseteq B$ and a bijection $f : A' \to B'$ such that:

Definition

Let $B = (I_{\alpha})_{\alpha \in A}$ and $C = (J_{\beta})_{\beta \in B}$ be barcodes.

A δ -matching between them consists of subsets $A' \subseteq A$, $B' \subseteq B$ and a bijection $f: A' \rightarrow B'$ such that:

- ▶ If $\alpha \notin A'$, $\beta \notin B'$, then length (I_{α}) , length $(J_{\beta}) < \delta$.
- ▶ If $f(I_{\alpha}) = J_{\beta}$, then the endpoints of I_{α} and J_{β} are within δ of eachother.

Definition

Let $B = (I_{\alpha})_{\alpha \in A}$ and $C = (J_{\beta})_{\beta \in B}$ be barcodes.

A δ -matching between them consists of subsets $A' \subseteq A$, $B' \subseteq B$ and a bijection $f: A' \rightarrow B'$ such that:

▶ If $\alpha \notin A'$, $\beta \notin B'$, then length(I_{α}), length(J_{β}) < δ .

• If $f(I_{\alpha}) = J_{\beta}$, then the endpoints of I_{α} and J_{β} are within δ of eachother. We define the *bottleneck distance* as

 $d_b(B, C) = \inf\{\delta > 0 \mid \text{there exists a } \delta\text{-matching between } B \text{ and } C\}.$

Definition

Let $B = (I_{\alpha})_{\alpha \in A}$ and $C = (J_{\beta})_{\beta \in B}$ be barcodes.

A δ -matching between them consists of subsets $A' \subseteq A$, $B' \subseteq B$ and a bijection $f: A' \rightarrow B'$ such that:

▶ If $\alpha \notin A'$, $\beta \notin B'$, then length (I_{α}) , length $(J_{\beta}) < 2\delta$.

• If $f(I_{\alpha}) = J_{\beta}$, then the endpoints of I_{α} and J_{β} are within δ of eachother. We define the *bottleneck distance* as

 $d_b(B, C) = \inf\{\delta > 0 \mid \text{there exists a } \delta\text{-matching between } B \text{ and } C\}.$



Theorem (Cohen-Steiner et al.)

Let $P, Q \subseteq \mathbb{R}^k$ be finite subsets and let B(P) and B(Q) be the barcodes of the persistent homology of their Rips filtrations. Then

 $d_b(B(P), B(Q)) \leq d_H(P, Q).$

Persistence diagrams

Definition If $B = (I_{\alpha})_{\alpha \in A}$ is a barcode, we define its *persistence diagram* as

$$\mathsf{dgm}(B) = ((\mathsf{inf}\ \mathit{I}_lpha, \mathsf{sup}\ \mathit{I}_lpha))_{lpha \in \mathcal{A}} \subseteq \mathbb{R}^2.$$

Persistence diagrams

Definition If $B = (I_{\alpha})_{\alpha \in A}$ is a barcode, we define its *persistence diagram* as

$$\mathsf{dgm}(B) = ((\mathsf{inf}\ \mathit{I}_lpha, \mathsf{sup}\ \mathit{I}_lpha))_{lpha \in \mathcal{A}} \subseteq \mathbb{R}^2.$$



Maximilian Schmahl

Contents

Homology of Spaces

Homology of Data

Ripser

Applications

C++ library by Ulrich Bauer to compute barcodes of Rips filtrations

C++ library by Ulrich Bauer to compute barcodes of Rips filtrations Also implemented in scikit-tda:

from ripser import ripser
from persim import plot_diagrams

Compute the persistence diagram of a Rips filtration # data is numpy array of points in euclidean space or a distance matrix dgm = ripser(data, maxdim = 1, thresh = inf, distance_matrix = False)

Plot the persistence diagram
plot_diagrams(dgm, show = False)
Contents

Homology of Spaces

Homology of Data

Ripser

Applications

 Infer something about the shape of a data set (e.g. Cosmic Microwave Background, Edelsbrunner et al.)

- Infer something about the shape of a data set (e.g. Cosmic Microwave Background, Edelsbrunner et al.)
- Infer something about the complexity of a data set (e.g. lung disease detection, Brodzki et al.)

- Infer something about the shape of a data set (e.g. Cosmic Microwave Background, Edelsbrunner et al.)
- Infer something about the complexity of a data set (e.g. lung disease detection, Brodzki et al.)
- Use it as an additional layer in machine learning methods

