IPAM, Los Angeles - April 30, 2019 Geometry of Big Data

The density of expected persistence diagrams and its kernel based estimation

Frédéric Chazal (Joint work with Vincent Divol)



https ://team.inria.fr/datashape/

https://geometrica.saclay.inria.fr/team/Fred.Chazal/

What is topological structure of data?





Challenges :

- \rightarrow no direct access to topological/geometric information : need of intermediate constructions (simplicial complexes);
- \rightarrow distinguish topological "signal" from noise;
- \rightarrow topological information may be multiscale;
- $\rightarrow\,$ statistical analysis of topological information.



What is topological structure of data?



Challenges :

- \rightarrow no direct access to topological/geometric information : need of intermediate constructions (simplicial complexes);
- \rightarrow distinguish topological "signal" from noise;
- \rightarrow topological information may be multiscale;
- $\rightarrow\,$ statistical analysis of topological information.

Topological Data Analysis (TDA) Persistent homology !

The classical TDA pipeline

Representations of persistence

- A filtered simplicial complex (or a filtration) S built on top of a set X is a family (S_a | a ∈ R) of subcomplexes of some fixed simplicial complex S with vertex set X s. t. S_a ⊆ S_b for any a ≤ b.
- More generaly, filtration = nested family of spaces.

Filtrations of simplicial complexes

- A filtered simplicial complex (or a filtration) S built on top of a set X is a family (S_a | a ∈ R) of subcomplexes of some fixed simplicial complex S with vertex set X s. t. S_a ⊆ S_b for any a ≤ b.
- More generaly, filtration = nested family of spaces.

Example : Let $(\mathbb{X}, d_{\mathbb{X}})$ be a metric space.

 The Vietoris-Rips filtration is the filtered simplicial complexe defined by : for a ∈ R,

 $[x_0, x_1, \cdots, x_k] \in \operatorname{Rips}(\mathbb{X}, a) \Leftrightarrow d_{\mathbb{X}}(x_i, x_j) \leq a, \text{ for all } i, j.$

X : metric data set

- Filtrations allow to construct "shapes" representing the data in a multiscale way.
- Persistent homology : encode the evolution of the topology across the scales → multi-scale topological signatures.
- Persistence diagrams are stable (w.r.t. Hausdorff metric).

X : metric data set

 $\operatorname{Filt}(\mathbb{X})$: filtered simplicial complex

- Filtrations allow to construct "shapes" representing the data in a multiscale way.
- Persistent homology : encode the evolution of the topology across the scales \rightarrow multi-scale topological signatures.
- Persistence diagrams are stable (w.r.t. Hausdorff metric).

X : metric data set

 $\operatorname{Filt}(\mathbb{X})$: filtered simplicial complex

- Filtrations allow to construct "shapes" representing the data in a multiscale way.
- Persistent homology : encode the evolution of the topology across the scales \rightarrow multi-scale topological signatures.
- Persistence diagrams are stable (w.r.t. Hausdorff metric).

► Filt(X) : filtered simplicial complex

- Filtrations allow to construct "shapes" representing the data in a multiscale way.
- Persistent homology : encode the evolution of the topology across the scales → multi-scale topological signatures.
- Persistence diagrams are stable (w.r.t. Hausdorff metric).

Persistence diagram

- the topology across the scales \rightarrow multi-scale topological signatures.
- Persistence diagrams are stable (w.r.t. Hausdorff metric).

Let $\mathbb{S} = (\mathbb{S}_a \mid a \in \mathbf{R})$ be a finite filtered simplicial complex with N simplices and let $\mathbb{S}_{a_1} \subset \mathbb{S}_{a_2} \subset \cdots \subset \mathbb{S}_{a_N}$ be the discrete filtration induced by the entering times of the simplices : $\mathbb{S}_{a_i} \setminus \mathbb{S}_{a_{i-1}} = \sigma_{a_i}$.

Let $\mathbb{S} = (\mathbb{S}_a \mid a \in \mathbf{R})$ be a finite filtered simplicial complex with N simplices and let $\mathbb{S}_{a_1} \subset \mathbb{S}_{a_2} \subset \cdots \subset \mathbb{S}_{a_N}$ be the discrete filtration induced by the entering times of the simplices : $\mathbb{S}_{a_i} \setminus \mathbb{S}_{a_{i-1}} = \sigma_{a_i}$.

Process the simplices according to their order of entrance in the filtration :

Let $k = \dim \sigma_{a_i}$

Let $\mathbb{S} = (\mathbb{S}_a \mid a \in \mathbf{R})$ be a finite filtered simplicial complex with N simplices and let $\mathbb{S}_{a_1} \subset \mathbb{S}_{a_2} \subset \cdots \subset \mathbb{S}_{a_N}$ be the discrete filtration induced by the entering times of the simplices : $\mathbb{S}_{a_i} \setminus \mathbb{S}_{a_{i-1}} = \sigma_{a_i}$.

Process the simplices according to their order of entrance in the filtration :

Let $k = \dim \sigma_{a_i}$

Case 1 : adding σ_{a_i} to $\mathbb{S}_{a_{i-1}}$ creates a new k-dimensional topological feature in \mathbb{S}_{a_i} (new homology class in H_k).

 \Rightarrow the birth of a k-dim feature is registered.

Let $\mathbb{S} = (\mathbb{S}_a \mid a \in \mathbf{R})$ be a finite filtered simplicial complex with N simplices and let $\mathbb{S}_{a_1} \subset \mathbb{S}_{a_2} \subset \cdots \subset \mathbb{S}_{a_N}$ be the discrete filtration induced by the entering times of the simplices : $\mathbb{S}_{a_i} \setminus \mathbb{S}_{a_{i-1}} = \sigma_{a_i}$.

Process the simplices according to their order of entrance in the filtration :

Let $k = \dim \sigma_{a_i}$

Case 1 : adding σ_{a_i} to $\mathbb{S}_{a_{i-1}}$ creates a new k-dimensional topological feature in \mathbb{S}_{a_i} (new homology class in H_k).

 \Rightarrow the birth of a k-dim feature is registered.

Case 2 : adding σ_{a_i} to $\mathbb{S}_{a_{i-1}}$ kills a (k-1)-dimensional topological feature in \mathbb{S}_{a_i} (homology class in H_{k-1}).

 \Rightarrow persistence algo. pairs the simplex σ_{a_i} to the simplex σ_{a_j} that gave birth to the killed feature.

Process the simplices according to their order of entrance in the filtration :

Let
$$k = \dim \sigma_{a_i}$$
 (ie. $\sigma_{a_i} = [v_0, \cdots, v_k]$)

Case 1 : adding σ_{a_i} to $\mathbb{S}_{a_{i-1}}$ creates a new k-dimensional topological feature in \mathbb{S}_{a_i} (new homology class in H_k).

 \Rightarrow the birth of a k-dim feature is registered.

Case 2 : adding σ_{a_i} to $\mathbb{S}_{a_{i-1}}$ kills a (k-1)-dimensional topological feature in \mathbb{S}_{a_i} (homology class in H_{k-1}).

 \Rightarrow persistence algo. pairs the simplex σ_{a_i} to the simplex σ_{a_j} that gave birth to the killed feature.

 $\rightarrow (\sigma_{a_j}, \sigma_{a_i})$: persistence pair

 \rightarrow $(a_j, a_i) \in \mathbb{R}^2$: point in the persistence diagram

Process the simplices according to their order of entrance in the filtration :

Let
$$k = \dim \sigma_{a_i}$$
 (ie. $\sigma_{a_i} = [v_0, \cdots, v_k]$)

Case 1 : adding σ_{a_i} to $\mathbb{S}_{a_{i-1}}$ creates a new k-dimensional topological feature in \mathbb{S}_{a_i} (new homology class in H_k).

 \Rightarrow the birth of a k-dim feature is registered.

Important to remember : the persistence pairs are determined by the order on the simplices; the corresponding $\rightarrow (a_i, a_i) \in \mathbb{R}^2$: point in the points in the diagrams are determined by the filtration indices.

Case 2 : adding σ_{a_i} to $\mathbb{S}_{a_{i-1}}$ kills a (k-1)-dimensional topological feature in \mathbb{S}_{a_i} (homology class in H_{k-1}).

 \Rightarrow persistence algo. pairs the simplex σ_{a_i} to the simplex σ_{a_i} that gave birth to the killed feature.

 $\rightarrow (\sigma_{a_i}, \sigma_{a_i})$: persistence pair

persistence diagram

Statistical setting

Statistical setting

What can be said about the distribution of diagrams $D[\mathcal{K}(\mathbb{X})]$? Understand the structure of $E[D[\mathcal{K}(\mathbb{X})]]$ in the non asymptotic setting ($|\mathbb{X}| = n$ is fixed, or bounded)

What does this mean?

Persistence diagrams as discrete measures

Motivations :

- The space of measures is much nicer that the space of P. D. !
- In the "standard" algebraic persistence theory, persistence diagrams naturally appear as discrete measures in the plane (over rectangles).
 [Chazal, de Silva, Glisse, Oudot 16]
- Many persistence representations can be expressed as

$$D(\phi) = \sum_{\mathbf{r} \in D} \phi(\mathbf{r}) = \int \phi(\mathbf{r}) dD(\mathbf{r})$$

for well-chosen functions ϕ .

Representation of Persistence diagrams

A representation is called linear if there exists $\phi: \mathbb{R}^2_> \to \mathcal{H}$ such that

$$\Phi(D) = \sum_{\mathbf{r} \in D} \phi(r) := D(\phi) = \int \phi(\mathbf{r}) \ dD(\mathbf{r})$$

In ML settings, well-suited linear representations of PD can be learnt.

[Hofer et al., NeurIPS 2017, Carrière et al, 2019]

Representation of Persistence diagrams

- D is a random persistence diagram
- E[D] is a deterministic measure on $\mathbb{R}^2_>$ defined by

$$\forall A \subset \mathbb{R}^2_>, \ E[D](A) = E[D(A)].$$

 $- D_1, \ldots, D_N$ i.i.d.

$$\overline{\Phi} = \frac{\Phi(D_1) + \dots + \Phi(D_N)}{N}$$
$$= \overline{\mu}(\phi)$$
$$\approx E[D](\phi)$$
$$E[D](\phi) = \int_{\mathbb{R}^2_{>}} \phi(\mathbf{r}) p(\mathbf{r}) d\mathbf{r}$$
?

Does E[D] has a density w.r.t. Lebesgue measure in \mathbb{R}^2 ? Estimation of p?

The density of expected persistence diagrams

Theorem : Fix $n \ge 1$. Assume that :

- *M* is a real analytic compact *d*-dimensional connected riemannian manifold possibly with boundary,
- X is a random variable on M^n having a density with respect to the Haussdorf measure \mathcal{H}_{dn} ,
- \mathcal{K} is the Vietoris-Rips filtration.

Then, for $s \ge 1$, $E[D_s[\mathcal{K}(\mathbb{X})]]$ has a density with respect to the Lebesgue measure on $\mathbb{R}^2_>$. Moreover, $E[D_0[\mathcal{K}(\mathbb{X})]]$ has a density with respect to the Lebesgue measure on the vertical line $\{0\} \times [0, \infty)$.

The density of expected persistence diagrams

Theorem : Fix $n \ge 1$. Assume that :

- *M* is a real analytic compact *d*-dimensional connected riemannian manifold possibly with boundary,
- X is a random variable on M^n having a density with respect to the Haussdorf measure \mathcal{H}_{dn} ,
- \mathcal{K} is the Vietoris-Rips filtration.

Then, for $s \ge 1$, $E[D_s[\mathcal{K}(\mathbb{X})]]$ has a density with respect to the Lebesgue measure on $\mathbb{R}^2_>$. Moreover, $E[D_0[\mathcal{K}(\mathbb{X})]]$ has a density with respect to the Lebesgue measure on the vertical line $\{0\} \times [0, \infty)$.

Theorem [smoothness]: Under the assumption of previous theorem, if moreover $\mathbb{X} \in M^n$ has a density of class C^k with respect to \mathcal{H}_{nd} . Then, for $s \geq 0$, the density of $E[D_s[\mathcal{K}(\mathbb{X})]]$ is of class C^k .

Remark : This is a particular case of a much more general result.

The density of expected persistence diagrams

Idea of the proof :

- Standard arguments from real analytic geometry : up to a set of measure 0, M^n can be decomposed into a finite set of open sets V_i on which the order on the simplices induced by the Rip filtration is constant.
- Classical argument from geometric measure theory (co-area formula): the map from V_i to the space of PD has maximal rank and the image of the random variable X has density with respect to Lebesgue measure on ℝ².

Filtrations revisited

Let n > 0 be an integer, \mathcal{F}_n : the collection of non-empty subsets of $\{1, \ldots, n\}$, M: a real analytic compact d-dim. connected manifold (poss. with boundary).

Filtering function :

$$\varphi = (\varphi[J])_{J \in \mathcal{F}_n} : M^n \to \mathbb{R}^{|\mathcal{F}_n|}$$

satisfiying the following conditions :

(K2) Invariance by permutation : For $J \in \mathcal{F}_n$ and for $(x_1, \ldots, x_n) \in M^n$, if τ is a permutation of the entries having support included in J, then $\varphi[J](x_{\tau(1)}, \ldots, x_{\tau(n)}) = \varphi[J](x_1, \ldots, x_n).$

(K3) Monotony : For $J \subset J' \in \mathcal{F}_n$, $\varphi[J] \leq \varphi[J']$.

Given $x = (x_1, \dots, x_n)$, $\varphi(x)$ induces an order on the faces of the simplex with n vertices that is a filtration $\mathcal{K}(x)$:

$$\forall J \in \mathcal{F}_n, \ J \in \mathcal{K}(x,r) \Longleftrightarrow \varphi[J](x) \le r.$$

The case of the Vietoris-Rips filtration

 $\varphi[J](x) = \max_{i,j \in J} d(x_i, x_j)$

- (K1) Absence of interaction : For $J \in \mathcal{F}_n$, $\varphi[J](x)$ only depends on x(J).
- (K2) Invariance by permutation : For $J \in \mathcal{F}_n$ and for $(x_1, \ldots, x_n) \in M^n$, if τ is a permutation of the entries having support included in J, then $\varphi[J](x_{\tau(1)}, \ldots, x_{\tau(n)}) = \varphi[J](x_1, \ldots, x_n).$
- (K3) Monotony : For $J \subset J' \in \mathcal{F}_n$, $\varphi[J] \leq \varphi[J']$.
- (K4) Compatibility : For a simplex $J \in \mathcal{F}_n$ and for $j \in J$, if $\varphi[J](x_1, \ldots, x_n)$ is not a function of x_j on some open set U of M^n , then $\varphi[J] \equiv \varphi[J \setminus \{j\}]$ on U.
- (K5') Smoothness : The function φ is subanalytic and the gradient of each of its entries J of size larger than 1 is non vanishing a.e. and for $J = \{j\}$, $\varphi[\{j\}] \equiv 0$.

Sketch of proof

1. There exists a partition of the complement of a (subanalytic) set of measure 0 in M^n by open sets V_1, \dots, V_R such that :

- the order of the simplices of $\mathcal{K}(x)$ is constant on each V_r ,
- for any $r=1,\cdots,R$, and any $x\in V_r$,

$$D_s[\mathcal{K}(x)] = \sum_{i=1}^{N_r} \delta_{\mathbf{r}_i}$$

with $\mathbf{r}_i = (\varphi[J_{i_1}](x), \varphi[J_{i_2}](x))$ where N_r , J_{i_1}, J_{i_2} only depends on V_r .

• J_{i_1}, J_{i_2} can be chosen so that the differential of

$$\Phi_{ir}: x \in V_r \to \mathbf{r}_i = (\varphi[J_{i_1}](x), \varphi[J_{i_2}](x))$$

has maximal rank 2.

Sketch of proof

2. The expected diagram can be written as

$$E[D_s[\mathcal{K}(\mathbb{X})]] = \sum_{r=1}^R E\left[\mathbb{1}\{\mathbb{X} \in V_r\} D_s[\mathcal{K}(\mathbb{X})]\right] = \sum_{r=1}^R E\left[\mathbb{1}\{\mathbb{X} \in V_r\} \sum_{i=1}^{N_r} \delta_{\mathbf{r}_i}\right]$$
$$= \sum_{r=1}^R \sum_{i=1}^{N_r} E\left[\mathbb{1}\{\mathbb{X} \in V_r\} \delta_{\mathbf{r}_i}\right]$$

Sketch of proof

2. The expected diagram can be written as

$$E[D_{s}[\mathcal{K}(\mathbb{X})]] = \sum_{r=1}^{R} E\left[\mathbb{1}\{\mathbb{X} \in V_{r}\}D_{s}[\mathcal{K}(\mathbb{X})]\right] = \sum_{r=1}^{R} E\left[\mathbb{1}\{\mathbb{X} \in V_{r}\}\sum_{i=1}^{N_{r}} \delta_{\mathbf{r}_{i}}\right]$$
$$= \sum_{r=1}^{R} \sum_{i=1}^{N_{r}} E\left[\mathbb{1}\{\mathbb{X} \in V_{r}\}\delta_{\mathbf{r}_{i}}\right]$$
$$\mu_{ir}$$
3. Use the co-area formula :
$$\mu_{ir}(B) = P(\Phi_{ir}(\mathbb{X}) \in B, \mathbb{X} \in V_{r})$$
$$= \int_{V_{r}} \mathbb{1}\{\Phi_{ir}(x) \in B\}\kappa(x)d\mathcal{H}_{nd}(x)$$
$$= \int_{U \in B} \int_{x \in \Phi_{ir}^{-1}(u)} (J\Phi_{ir}(x))^{-1}\kappa(x)d\mathcal{H}_{nd-2}(x)du.$$
Density of μ_{ir}

The Hausdorff measure and the co-area formula

Definition : Let k be a non-negative number. For $A \subset \mathbb{R}^D$, and $\delta > 0$, consider

$$\mathcal{H}_k^{\delta}(A) := \inf \left\{ \sum_i \operatorname{diam}(U_i)^k, A \subset \bigcup_i U_i \text{ and } \operatorname{diam}(U_i) < \delta \right\}.$$

The *k*-dimensional Haussdorf measure on \mathbb{R}^D of A is defined by $\mathcal{H}_k(A) := \lim_{\delta \to 0} \mathcal{H}_k^{\delta}(A)$.

Theorem [Co-area formula] : Let M (resp. N) be a smooth Riemannian manifold of dimension m (resp n). Assume that $m \ge n$ and let $\Phi : M \to N$ be a differentiable map. Denote by $D\Phi$ the differential of Φ . The Jacobian of Φ is defined by $J\Phi = \sqrt{\det((D\Phi) \times (D\Phi)^t)}$. For $f : M \to N$ a positive measurable function, the following equality holds :

$$\int_{M} f(x) J\Phi(x) d\mathcal{H}_{m}(x) = \int_{N} \left(\int_{x \in \Phi^{-1}(\{y\})} f(x) d\mathcal{H}_{m-n}(x) \right) d\mathcal{H}_{n}(y).$$

Persistence images

[Adams et al, JMLR 2017]

For $K : \mathbb{R}^2 \to \mathbb{R}$ a kernel and H a bandwidth matrix (e.g. a symmetric positive definite matrix), pose for $u \in \mathbb{R}^2$, $K_H(z) = |H|^{-1/2} K(H^{-1/2} \cdot u)$

For $D = \sum_i \delta_{\mathbf{r}_i}$ a diagram, $K : \mathbb{R}^2 \to \mathbb{R}$ a kernel, H a bandwidth matrix and $w : \mathbb{R}^2 \to \mathbb{R}_+$ a weight function, one defines the persistence surface of D with kernel K and weight function w by :

$$\forall z \in \mathbb{R}^2, \ \rho(D)(u) = \sum_i w(\mathbf{r}_i) K_H(u - \mathbf{r}_i) = D(wK_H(u - \cdot))$$

Persistence images

[Adams et al, JMLR 2017]

For $K : \mathbb{R}^2 \to \mathbb{R}$ a kernel and H a bandwidth matrix (e.g. a symmetric positive definite matrix), pose for $u \in \mathbb{R}^2$, $K_H(z) = |H|^{-1/2} K(H^{-1/2} \cdot u)$

For $D = \sum_i \delta_{\mathbf{r}_i}$ a diagram, $K : \mathbb{R}^2 \to \mathbb{R}$ a kernel, H a bandwidth matrix and $w : \mathbb{R}^2 \to \mathbb{R}_+$ a weight function, one defines the persistence surface of D with kernel K and weight function w by :

$$\forall z \in \mathbb{R}^2, \ \rho(D)(u) = \sum_i w(\mathbf{r}_i) K_H(u - \mathbf{r}_i) = D(wK_H(u - \cdot))$$

 \Rightarrow persistence surfaces can be seen as kernel based estimators of $E[D_s[\mathcal{K}(\mathbb{X})]]$.

Persistence images

The realization of 3 different processes

The overlay of 40 different persistence diagrams

The persistence images with weight function $w(\mathbf{r}) = (r_2 - r_1)^3$ and bandwith selected using cross-validation.

Thank you for your attention

References :

• F. Chazal, V. Divol, *The density of expected persistence diagrams and its kernel based estimation*, SoCG 2018.

Software :

- GUDHI library C++ / Python : http ://gudhi.gforge.inria.fr/
- R package TDA : Statistical Tools for Topological Data Analysis