Journées de Géométrie Algorithmique Aussois, Décembre 2017

A theoretical framework for the analysis of Reeb graph and Mapper

Mathieu Carrière — joint work with B. Michel and S. Oudot





UMR 6629 - Nantes

Reeb Graphs and Mapper



Reeb Graphs and Mapper



- $\longrightarrow \text{visualization}$
- $\longrightarrow \mathsf{clustering}$
- $\longrightarrow \text{feature selection}$



- \longrightarrow visualization
- $\longrightarrow \mathsf{clustering}$
- $\longrightarrow \text{feature selection}$

Principle: identify statistically relevant subpopulations through **patterns** (flares, loops)





breast cancer subtype identification

[Nicolau et al. 2011]



protein folding pathways

[Yao et al. 2009]



Data Skeletonization

[Ge et al. 2011]



Burning regions of lean hydrogen flame over time [Weber et al. 2011]

Reeb Graphs

 $x \sim y \iff [f(x) = f(y) \text{ and } x, y \text{ belong to same cc of } f^{-1}(f(x))]$

Def: $\operatorname{R}_f(X) = X/\sim$





Prop: $R_f(X)$ is a graph if (X, f) is of **Morse type**

 $x \sim y \iff [f(x) = f(y) \text{ and } x, y \text{ belong to same cc of } f^{-1}(f(x))]$ Def: $R_f(X) = X/\sim$

Caveat: computation from point cloud is difficult



 $x \sim y \iff [f(x) = f(y) \text{ and } x, y \text{ belong to same cc of } f^{-1}(f(x))]$ Def: $R_f(X) = X/\sim$











Input:

- topological space \boldsymbol{X}
- continuous function $f:X\to \mathbb{R}$
- cover ${\mathcal I}$ of $\operatorname{im}(f)$ by open intervals: $\operatorname{im}(f)\subseteq \bigcup_{I\in {\mathcal I}} I$

Method:

- Compute *pullback cover* \mathcal{U} of X: $\mathcal{U} = \{f^{-1}(I)\}_{I \in \mathcal{I}}$
- \bullet Refine ${\mathcal U}$ by separating the connected components
- The Mapper is the *nerve* of \mathcal{V} :
 - 1 vertex per element $V \in \mathcal{V}$
 - 1 edge per intersection $V \cap V' \neq \emptyset$, $V,V' \in \mathcal{V}$
 - 1 k-simplex per (k+1)-fold intersection $\bigcap_{i=0}^{k} V_i \neq \emptyset$, $V_0, \cdots, V_k \in \mathcal{V}$

Mapper (discrete setting)

Input:

- point cloud $P \subseteq X$ with metric d_P
- continuous function $f: \textbf{\textit{P}} \rightarrow \mathbb{R}$
- cover \mathcal{I} of $\operatorname{im}(f)$ by open intervals: $\operatorname{im} f \subseteq \bigcup_{I \in \mathcal{I}} I$

Method:

- Compute *pullback cover* \mathcal{U} of P: $\mathcal{U} = \{f^{-1}(I)\}_{I \in \mathcal{I}}$
- Refine ${\mathcal U}$ by separating each of its elements into its various connected components in $G\to$ connected cover ${\mathcal V}$
- The Mapper is the *nerve* of \mathcal{V} :

(intersections materialized by data points)

- 1 vertex per element $V \in \mathcal{V}$
- 1 edge per intersection $V \cap V' \neq \emptyset$, $V,V' \in \mathcal{V}$
- 1 k-simplex per (k+1)-fold intersection $\bigcap_{i=0}^{k} V_i \neq \emptyset$, $V_0, \cdots, V_k \in \mathcal{V}$

Mapper (discrete setting)









In practice: trial-and-error

 $\hat{f} = \mathsf{density} \ \mathsf{estimator}$

r = 0.3, g = 20%









 $\delta = 0.1\%$

 $\delta = 1\%$

 $\delta = 10\%$

Examples



Examples





- family of *excursion sets* (sublevel then superlevel sets) of Reeb graph
- use *homological algebra* to encode the evolution of the topology of the family



- family of *excursion sets* (sublevel then superlevel sets) of Reeb graph
- use *homological algebra* to encode the evolution of the topology of the family



- family of *excursion sets* (sublevel then superlevel sets) of Reeb graph
- use *homological algebra* to encode the evolution of the topology of the family



- family of *excursion sets* (sublevel then superlevel sets) of Reeb graph
- use *homological algebra* to encode the evolution of the topology of the family



- family of *excursion sets* (sublevel then superlevel sets) of Reeb graph
- use *homological algebra* to encode the evolution of the topology of the family



- family of *excursion sets* (sublevel then superlevel sets) of Reeb graph
- use *homological algebra* to encode the evolution of the topology of the family



- family of *excursion sets* (sublevel then superlevel sets) of Reeb graph
- use *homological algebra* to encode the evolution of the topology of the family



- family of *excursion sets* (sublevel then superlevel sets) of Reeb graph
- use *homological algebra* to encode the evolution of the topology of the family



- family of *excursion sets* (sublevel then superlevel sets) of Reeb graph
- use *homological algebra* to encode the evolution of the topology of the family



- family of *excursion sets* (sublevel then superlevel sets) of Reeb graph
- use *homological algebra* to encode the evolution of the topology of the family



- family of *excursion sets* (sublevel then superlevel sets) of Reeb graph
- use *homological algebra* to encode the evolution of the topology of the family


- family of *excursion sets* (sublevel then superlevel sets) of Reeb graph
- use *homological algebra* to encode the evolution of the topology of the family



- family of *excursion sets* (sublevel then superlevel sets) of Reeb graph
- use *homological algebra* to encode the evolution of the topology of the family



- family of *excursion sets* (sublevel then superlevel sets) of Reeb graph
- use *homological algebra* to encode the evolution of the topology of the family



- family of *excursion sets* (sublevel then superlevel sets) of Reeb graph
- use *homological algebra* to encode the evolution of the topology of the family



- family of *excursion sets* (sublevel then superlevel sets) of Reeb graph
- use *homological algebra* to encode the evolution of the topology of the family



- family of *excursion sets* (sublevel then superlevel sets) of Reeb graph
- use *homological algebra* to encode the evolution of the topology of the family



- family of *excursion sets* (sublevel then superlevel sets) of Reeb graph
- use *homological algebra* to encode the evolution of the topology of the family



- family of *excursion sets* (sublevel then superlevel sets) of Reeb graph
- use *homological algebra* to encode the evolution of the topology of the family



- family of *excursion sets* (sublevel then superlevel sets) of Reeb graph
- use *homological algebra* to encode the evolution of the topology of the family



Metric for Extended Persistence Diagrams



Metric for Extended Persistence Diagrams

Partial matching $M : Dg \leftrightarrow Dg'$ Matched pair $(x, y) \in M : c(x, y) = ||x - y||_{\infty}$ Unmatched point $z \in X \sqcup Y : c(z) = ||z - \overline{z}||_{\infty}$ $c(M) = \max\{\max_{(x, y)} c(x, y), \max_{z} c(z)\}$



Metric for Extended Persistence Diagrams

Partial matching $M : Dg \leftrightarrow Dg'$ Matched pair $(x, y) \in M : c(x, y) = ||x - y||_{\infty}$ Unmatched point $z \in X \sqcup Y : c(z) = ||z - \overline{z}||_{\infty}$ $c(M) = \max\{\max_{(x, y)} c(x, y), \max_{z} c(z)\}$





Metric Properties

Thm (stability): [Bauer, Ge, Wang 2013] $d_B(Dg(R_f), Dg(R_g)) \le 6 d_{GH}(R_f, R_g)$



Metric Properties

Thm (stability): [Bauer, Ge, Wang 2013] $d_B(Dg(R_f), Dg(R_g)) \le 6 d_{GH}(R_f, R_g)$



Thm (stability): [Bauer, Ge, Wang 2013] $d_B(Dg(R_f), Dg(R_g)) \le 6 d_{GH}(R_f, R_g)$





Thm (stability): [Bauer, Ge, Wang 2013] $d_B(Dg(R_f), Dg(R_g)) \le 6 d_{GH}(R_f, R_g)$

Note: $d_B(Dg(\cdot), Dg(\cdot))$ is only a pseudometric on Reeb graphs **Thm:** [C., Oudot 2016] $d_B(Dg(\cdot), Dg(\cdot))$ is *locally* a metric equivalent to d_{GH}



Reminder: Mapper \equiv *pixelized* Reeb graph







Def: Given X, f, \mathcal{I} : $\operatorname{Dg}(\operatorname{M}_{f}) = \left(\operatorname{DB}(\operatorname{R}_{f}) \setminus Q_{\mathcal{I}}^{\operatorname{DB}}\right) \cup \left(\operatorname{UB}(\operatorname{R}_{f}) \setminus Q_{\mathcal{I}}^{\operatorname{UB}}\right) \cup \left(\operatorname{L}(\operatorname{R}_{f}) \setminus Q_{\mathcal{I}}^{\operatorname{L}}\right)$

0

Def: Given X, f, \mathcal{I} : $\mathrm{Dg}(\mathrm{M}_f) = \left(\mathrm{DB}(\mathrm{R}_f) \setminus Q_{\mathcal{I}}^{\mathrm{DB}}\right) \cup \left(\mathrm{UB}(\mathrm{R}_f) \setminus Q_{\mathcal{I}}^{\mathrm{UB}}\right) \cup \left(\mathrm{L}(\mathrm{R}_f) \setminus Q_{\mathcal{I}}^{\mathrm{L}}\right)$

























Cor: $Dg(M_f) = Dg(R_f)$ whenever the resolution r of \mathcal{I} is smaller than the smallest distance from $Dg(R_f) \setminus \Delta$ to the diagonal Δ



Cor: $Dg(M_f) = Dg(R_f)$ whenever the resolution r of \mathcal{I} is smaller than the smallest distance from $Dg(R_f) \setminus \Delta$ to the diagonal Δ

Thm: [C., Oudot 2017] $d_{\mathrm{GH}}(\mathrm{M}_f(X,\mathcal{I}),\mathrm{R}_f(X)) \leq 3r$







Statistics for Mapper







Prop: [C., Michel, Oudot 2017] $\widehat{\mathrm{M}}_n = \mathrm{M}_f(\widehat{X}_n)$ is measurable



Prop: [C., Michel, Oudot 2017]
$$\widehat{\mathrm{M}}_n = \mathrm{M}_f(\widehat{X}_n)$$
 is measurable

Goal: Find heuristics to compute "good" δ_n , g_n , r_n

Assess quality through confidence regions and convergence rates





Confidence regions: given $\alpha \in (0,1)$, find $c_n(\alpha) \ge 0$ s.t.: $\lim_{n \to \infty} \mathbb{P}\left(d_B\left(\widehat{M}_n, R_f(X)\right) > c_n(\alpha)\right) \le \alpha$ $\to d_B\text{-ball of radius } c_n(\alpha) \text{ around } \mathrm{Dg}(\widehat{M}_n)$







Confidence regions: given $\alpha \in (0, 1)$, find $c_n(\alpha) \ge 0$ s.t.: $\limsup_{n \to \infty} \mathbb{P}\left(d_B\left(\widehat{M}_n, R_f(X)\right) > c_n(\alpha)\right) \le \alpha$ $\to d_B\text{-ball of radius } c_n(\alpha) \text{ around } \mathrm{Dg}\left(\widehat{M}_n\right)$

Convergence Rate: estimate $\mathbb{E}\left[d_B(\widehat{M}_n, R_f(X))\right]$ w.r.t. n



Regularity of the filter function: (exact) modulus of continuity of f $\omega(\delta) = \sup_{\|x-x'\| \le \delta} |f(x) - f(x')|$

Approximation inequality: [C., Michel, Oudot 2017] Let $\hat{X}_n \subset X$. Under some *regularity assumptions* on X, f, δ, r, g , one has:

$$d_B\left(R_f(X), M_f(\widehat{X}_n)\right) \le r + 2\omega(\delta)$$

Rate of Convergence



 $4d_{\mathrm{H}}(\widehat{X}_{n}, X) \leq \delta \leq C(X)$ $\max\{|f(X_{i}) - f(X_{j})| : ||X_{i} - X_{j}|| \leq \delta\} < gr$

Approximation inequality: [C., Michel, Oudot 2017] Let $\widehat{X}_n \subset X$. Under some *regularity assumptions* on X, f, δ, r, g , one has: $d_B\left(R_f(X), M_f(\widehat{X}_n)\right) \leq r + 2\omega(\delta)$

Rate of Convergence



 $4d_{H}(\widehat{X}_{n}, X) \leq \delta \leq C(X)$ $\max\{\max\{|f(X_{i}) - f(X_{j})|, |\widehat{f}(X_{i}) - \widehat{f}(X_{j})|\} : \|X_{i} - X_{j}\| \leq \delta\} \leq rg$

Approximation inequality: [C., Michel, Oudot 2017] Let $\widehat{X}_n \subset X$. Under some *regularity assumptions* on X, f, δ, r, g , one has: $d_B\left(R_f(X), M_{\widehat{f}}(\widehat{X}_n)\right) \leq 2r + 2\omega(\delta) + \max\{|f(X_i) - \widehat{f}(X_i)|\}$

Rate of Convergence $\overbrace{f}^{(X, d_X, \mu)} \quad n \text{ points sampled} \\ i.i.d. according to \mu$

 $V_n = \max\{f(X_i) - f(X_j) : \|X_i - X_j\| \le \delta_n\}$

Thm: [C., Michel, Oudot 2017] If μ is (a, b)-standard, then for $\delta_n = 4\left(\frac{2\log n}{an}\right)^{1/b}$, $g_n \in \left(\frac{1}{3}, \frac{1}{2}\right)$, $r_n = \frac{V_n}{g_n}$, one has: $\sup \mathbb{E}\left[d_{\mathbb{P}}\left(M_n(\widehat{X}_n) \operatorname{R}_n(X)\right)\right] \leq \omega \left(\frac{\log n}{2}\right)^{1/b}$

$$\sup_{u \in \mathcal{P}} \mathbb{E} \left[\mathrm{d}_B \left(\mathrm{M}_f(\widehat{X}_n), \mathrm{R}_f(X) \right) \right] \lesssim \omega \left(\frac{\log n}{n} \right)^{-1}$$

Rate of Convergence $\overbrace{f}^{(X, d_X, \mu)} \quad n \text{ points sampled} \\ i.i.d. according to \mu$

 $V_n = \max\{f(X_i) - f(X_j) : \|X_i - X_j\| \le \delta_n\}$

Thm: [C., Michel, Oudot 2017] If μ is (a, b)-standard, then for $\delta_n = 4\left(\frac{2\log n}{\varpi n}\right)^{1/b}$, $g_n \in \left(\frac{1}{3}, \frac{1}{2}\right)$, $r_n = \frac{V_n}{g_n}$, one has: $\sup_{\mu \in \mathcal{P}} \mathbb{E}\left[d_B\left(M_f(\widehat{X}_n), R_f(X)\right)\right] \lesssim \omega \left(\frac{\log n}{n}\right)^{1/b}$


Subsampling to tune δ_n : let $\beta > 0$ and take $s(n) = n \log(n)^{-(1+\beta)}$ $\delta_n = d_H(\hat{X}_n^{s(n)}, \hat{X}_n)$ where $\hat{X}_n^{s(n)} \subset \hat{X}_n$ of size s(n)

Rate of Convergence \widehat{X}_{n} \widehat{X}_{n} $\widehat{\delta}_{n}$ \widehat{X}_{n} $\widehat{\delta}_{n}$ \widehat{g}_{n}, r_{n} $\widehat{M}_{f}(\widehat{X}_{n})$

Subsampling to tune δ_n : let $\beta > 0$ and take $s(n) = n \log(n)^{-(1+\beta)}$ $\delta_n = d_H(\hat{X}_n^{s(n)}, \hat{X}_n)$ where $\hat{X}_n^{s(n)} \subset \hat{X}_n$ of size s(n)

Thm: [C., Michel, Oudot 2017] If μ is (a, b)-standard, then for δ_n , $g_n \in \left(\frac{1}{3}, \frac{1}{2}\right)$, $r_n = \frac{V_n}{g_n}$, one has $\sup_{\mu \in \mathcal{P}} \mathbb{E} \left[d_B \left(M_f(\widehat{X}_n), R_f(X) \right) \right] \lesssim \omega \left(\frac{\log(n)^{2+\beta}}{n} \right)^{1/b}$

Rate of Convergence $\overbrace{f}^{(X, d_X, \mu)} \quad n \text{ points sampled} \\ \overbrace{i.i.d. according to } \mu \quad \overbrace{f}^{(X, d_X, \mu)} \quad f \xrightarrow{f}^{(X, d_X, \mu)}$

Subsampling to tune δ_n : let $\beta > 0$ and take $s(n) = n \log(n)^{-(1+\beta)}$ $\delta_n = d_H(\hat{X}_n^{s(n)}, \hat{X}_n)$ where $\hat{X}_n^{s(n)} \subset \hat{X}_n$ of size s(n)

Thm: [C., Michel, Oudot 2017] If μ is (a, b)-standard, then for δ_n , $g_n \in \left(\frac{1}{3}, \frac{1}{2}\right)$, $r_n = \frac{\max\{V_n, \widehat{V}_n\}}{g_n}$, one has $\sup_{\mu \in \mathcal{P}} \mathbb{E}\left[d_B\left(M_{\widehat{f}}(\widehat{X}_n), R_f(X)\right)\right] \lesssim \omega \left(\frac{\log(n)^{2+\beta}}{n}\right)^{1/b} + \mathbb{E}\left[\max|f(X) - \widehat{f}(X)|\right]$



Minimax Optimality: [C., Michel, Oudot 2017] for any estimator $\widehat{\mathrm{R}}_n$,

$$\omega\left(\frac{1}{n}\right)^{1/b} \lesssim \sup_{\mu \in \mathcal{P}} \mathbb{E}\left[d_B\left(\widehat{\mathbf{R}}_n, \mathbf{R}_f(X)\right)\right]$$

Rate of Convergence



Ex : PCA filter

 Π_1 : 1st principal direction of covariance operator $\widehat{\Pi}_1$: 1st principal direction of empirical covariance operator Using [Biau et. al. 2012]:

$$\mathbb{E}\left[\mathrm{d}_B\left(\mathrm{R}_{\Pi_1}(X), \mathrm{M}_{\widehat{\Pi}_1}(\widehat{X}_n)\right)\right] \lesssim \left(\frac{\log(n)^{2+\beta}}{n}\right)^{1/b} \vee \frac{1}{\sqrt{n}}$$

Confidence Regions

Either from proof of previous result with:

$$\mathbb{E}\left[\mathrm{d}_B\left(\mathrm{M}_f(\widehat{X}_n), \mathrm{R}_f(X)\right)\right] = \int_{\alpha} \mathbb{P}\left(\mathrm{d}_B\left(\mathrm{M}_f(\widehat{X}_n), \mathrm{R}_f(X)\right) \ge \alpha\right) \mathrm{d}\alpha$$

Or bootstrap (only empirical):

- draw $\widehat{X}_n^* = X_1^*, \cdots, X_n^*$ iid from $\mu_{\widehat{X}_n}$ (empirical measure on \widehat{X}_n)
- compute $d^* = d_B\left(M_f(\widehat{X}_n^*), M_f(\widehat{X}_n)\right)$
- repeat N times to get $\mathbf{d}_1^*,\cdots,\mathbf{d}_N^*$
- let q_{α} be the (1α) quantile of $\frac{1}{N} \sum_{i=1}^{N} \mathbb{I}(\sqrt{n} d_i^* \ge t)$

• take
$$c_n(\alpha) = \frac{q_\alpha}{\sqrt{n}}$$

Experiments



Experiments









































Conclusion

Structure and Stability of Mapper

Parameter Selection for Mapper

Extensions:

Multivariate function $f: X \to \mathbb{R}^n$

Space of Reeb graphs (curvature, barycenters, interpolation...)

Conclusion

Structure and Stability of Mapper

Parameter Selection for Mapper

Extensions:

Multivariate function $f: X \to \mathbb{R}^n$

Space of Reeb graphs (curvature, barycenters, interpolation...)

Thank you!