



INTRINSIC PERSISTENT HOMOLOGY VIA DENSITY-BASED METRIC LEARNING

XIMENA FERNÁNDEZ*

joint work with E. Borghini, P. Groisman and G. Mindlin

38TH WORKSHOP IN GEOMETRIC TOPOLOGY

16th June 2021

*EPSRC Centre for Topological Data Analysis

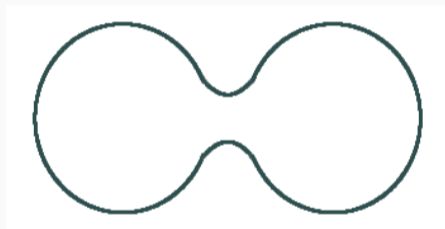


Prifysgol
Abertawe
Swansea
University

The problem

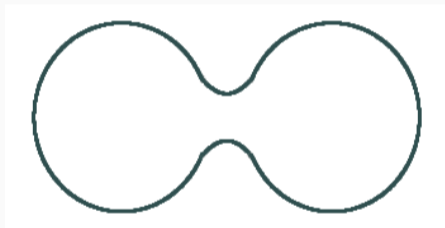
Homology inference

(\mathcal{M}, g) a d -dimensional Riemannian manifold
embedded in \mathbb{R}^D .



Homology inference

(\mathcal{M}, g) a d -dimensional Riemannian manifold embedded in \mathbb{R}^D .

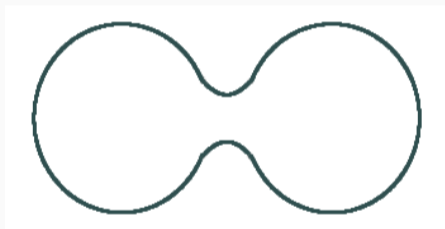


$\mathbb{X}_n = \{x_1, x_2, \dots, x_n\}$ a finite sample of \mathcal{M} .



Homology inference

(\mathcal{M}, g) a d -dimensional Riemannian manifold embedded in \mathbb{R}^D .



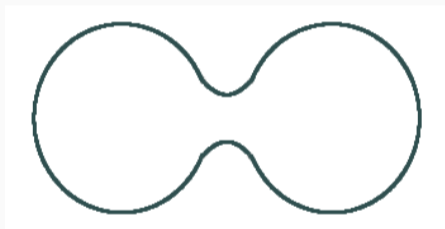
Q: How to infer the homology of \mathcal{M} from the sample \mathbb{X}_n ?

$\mathbb{X}_n = \{x_1, x_2, \dots, x_n\}$ a finite sample of \mathcal{M} .



Homology inference

(\mathcal{M}, g) a d -dimensional Riemannian manifold embedded in \mathbb{R}^D .



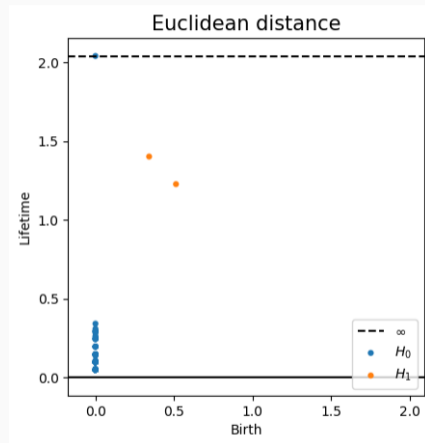
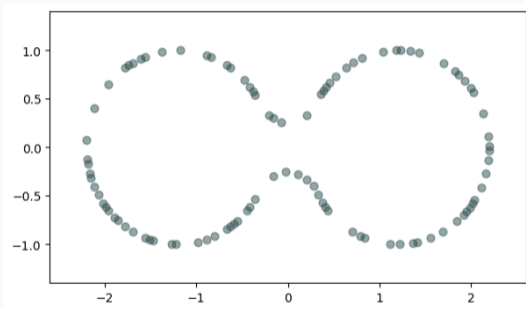
Q: How to infer the homology of \mathcal{M} from the sample \mathbb{X}_n ?

$\mathbb{X}_n = \{x_1, x_2, \dots, x_n\}$ a finite sample of \mathcal{M} .

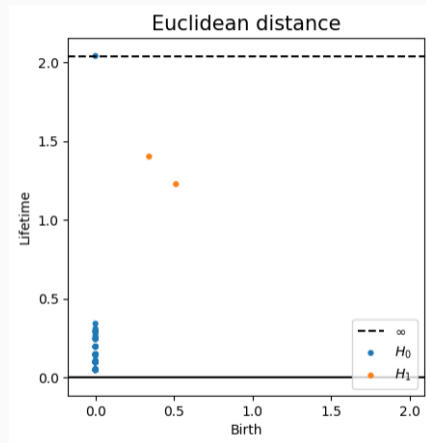
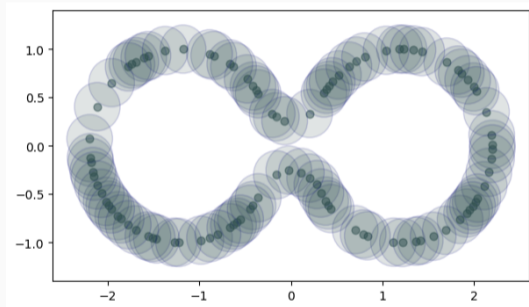


A: Compute persistent homology of \mathbb{X}_n .

Ambient persistent homology

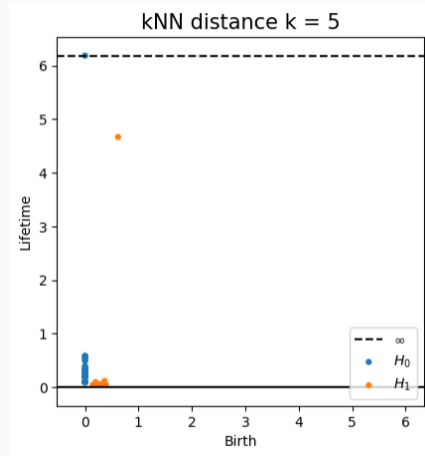
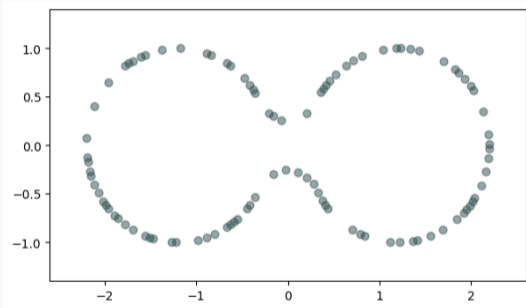


Ambient persistent homology

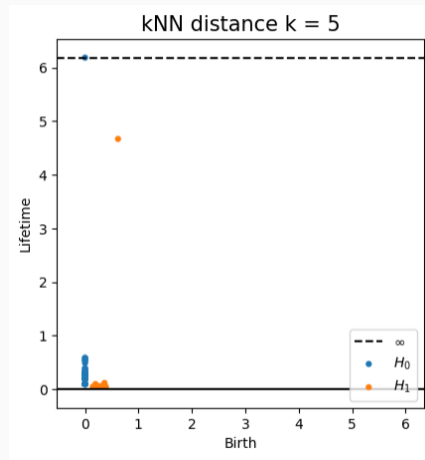
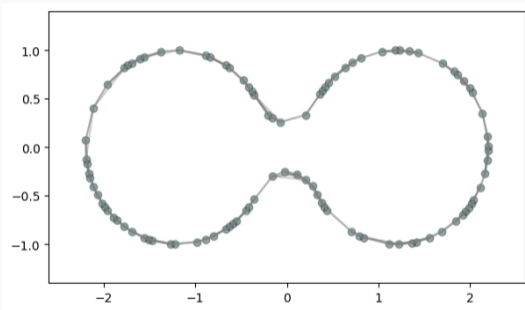


- $\text{Rips}_\epsilon(\mathcal{M}, d_E) \simeq \mathcal{M}$ for $\epsilon < 2\text{rch}(\mathcal{M})$

Intrinsic persistent homology

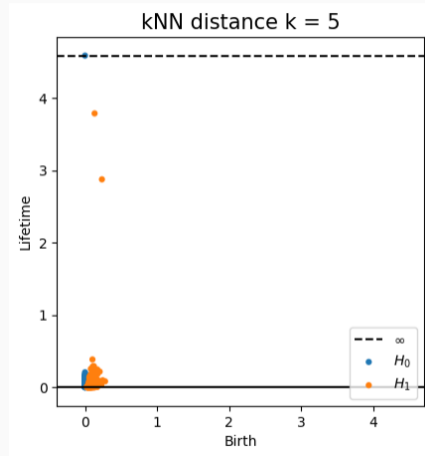
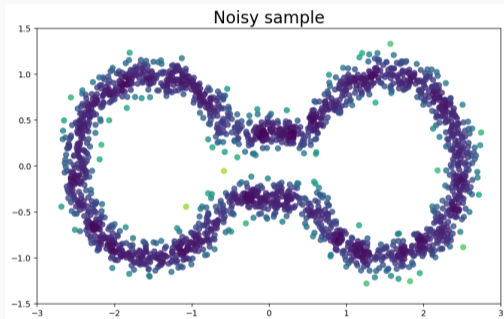


Intrinsic persistent homology

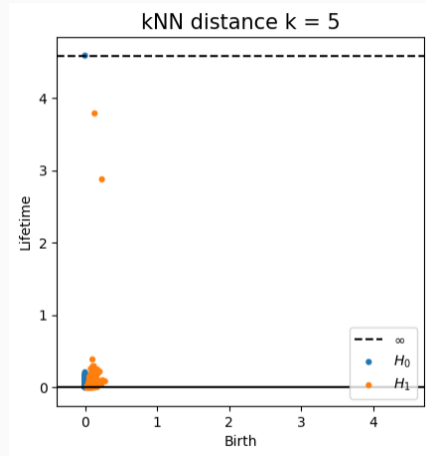
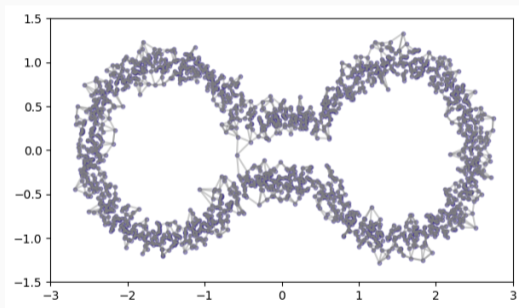


- $\text{Rips}_\epsilon(\mathcal{M}, d_{\mathcal{M}}) \simeq \mathcal{M}$ for $\epsilon < \text{conv}(\mathcal{M})$

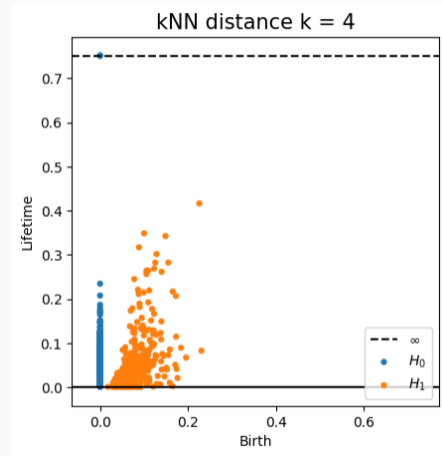
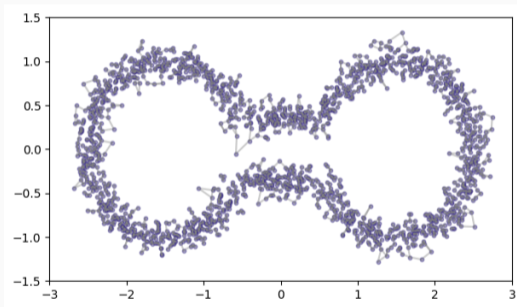
The problem of noise



The problem of noise



The problem of noise



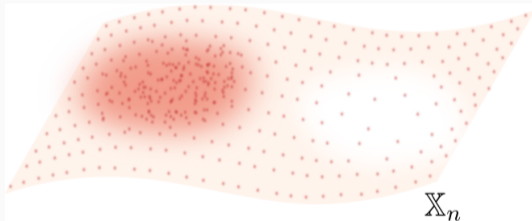
Density-based manifold learning

The manifold (and density) assumption

$\mathbb{X}_n = \{x_1, x_2, \dots, x_n\}$ a finite set of points in \mathbb{R}^D .

We assume that:

- * \mathbb{X}_n lies in a d -dimensional Riemannian manifold \mathcal{M} embedded in \mathbb{R}^D ,
- ** \mathbb{X}_n is drawn according to a smooth density $f : \mathcal{M} \rightarrow \mathbb{R}_{>0}$.

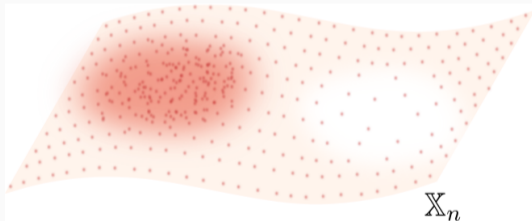


The manifold (and density) assumption

$\mathbb{X}_n = \{x_1, x_2, \dots, x_n\}$ a finite set of points in \mathbb{R}^D .

We assume that:

- * \mathbb{X}_n lies in a d -dimensional Riemannian manifold \mathcal{M} embedded in \mathbb{R}^D ,
- ** \mathbb{X}_n is drawn according to a smooth density $f : \mathcal{M} \rightarrow \mathbb{R}_{>0}$.



Idea:

- Consider a Riemannian metric that depends on f .
- Find an estimator of the (density-based) Riemannian metric from the sample.

Deformed Riemannian metric

- Let (\mathcal{M}, g) be a Riemannian manifold and let $f : \mathcal{M} \rightarrow \mathbb{R}_{>0}$ be a smooth density.

Deformed Riemannian metric

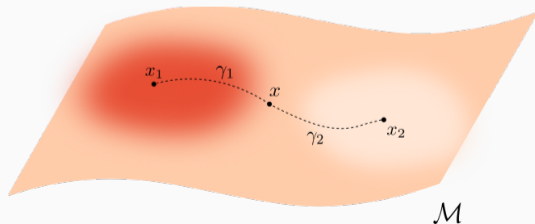
- Let (\mathcal{M}, g) be a Riemannian manifold and let $f : \mathcal{M} \rightarrow \mathbb{R}_{>0}$ be a smooth density.
- For $q > 0$, and consider the deformed metric tensor $g_q = f^{-2q}g$.

Deformed Riemannian metric

- Let (\mathcal{M}, g) be a **Riemannian manifold** and let $f : \mathcal{M} \rightarrow \mathbb{R}_{>0}$ be a smooth **density**.
- For $q > 0$, and consider the **deformed metric tensor** $g_q = f^{-2q}g$.
- The induced **deformed Riemannian distance** in \mathcal{M} is

$$d_{f,q}(x, y) = \inf_{\gamma} \int_I \frac{1}{f(\gamma_t)^q} \sqrt{g(\dot{\gamma}_t, \dot{\gamma}_t)} dt$$

over all $\gamma : I \rightarrow \mathcal{M}$ with $\gamma(0) = x$ and $\gamma(1) = y$.

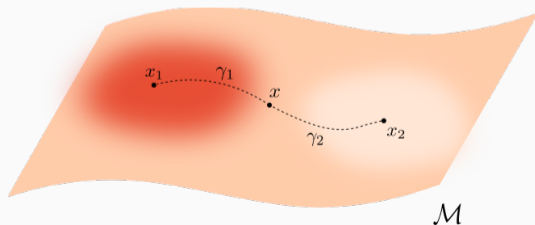


Deformed Riemannian metric

- Let (\mathcal{M}, g) be a **Riemannian manifold** and let $f : \mathcal{M} \rightarrow \mathbb{R}_{>0}$ be a smooth **density**.
- For $q > 0$, and consider the **deformed metric tensor** $g_q = f^{-2q}g$.
- The induced **deformed Riemannian distance** in \mathcal{M} is

$$d_{f,q}(x, y) = \inf_{\gamma} \int_I \frac{1}{f(\gamma_t)^q} \sqrt{g(\dot{\gamma}_t, \dot{\gamma}_t)} dt$$

over all $\gamma : I \rightarrow \mathcal{M}$ with $\gamma(0) = x$ and $\gamma(1) = y$.



Fermat distance

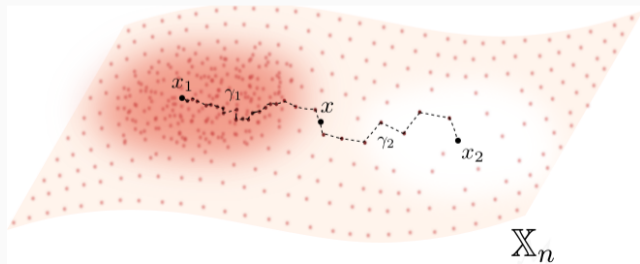
- Let $\mathbb{X}_n \subseteq \mathbb{R}^D$ a **sample** of points.

Fermat distance

- Let $\mathbb{X}_n \subseteq \mathbb{R}^D$ a **sample** of points.
- For $p > 1$, the **(sample) Fermat distance** between $x, y \in \mathbb{R}^D$ is defined by

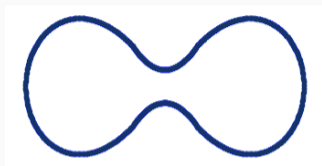
$$d_{\mathbb{X}_n, p}(x, y) = \inf_{\gamma} \sum_{i=0}^r |x_{i+1} - x_i|^p$$

over all paths $\gamma = (x_0, \dots, x_{r+1})$ of finite length with $x_0 = x$, $x_{r+1} = y$ and $\{x_1, x_2, \dots, x_r\} \subseteq \mathbb{X}_n$.



Example (Fermat distance)

Manifold

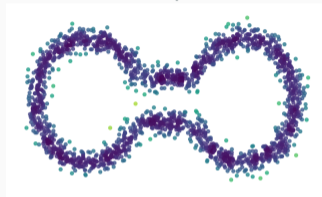


Example (Fermat distance)

Manifold

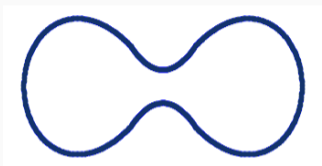


Sample

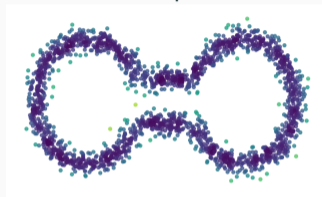


Example (Fermat distance)

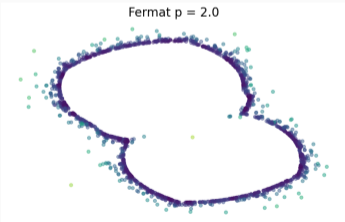
Manifold



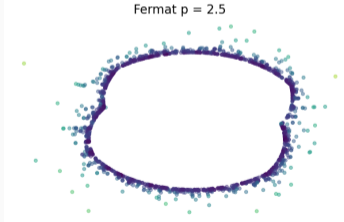
Sample



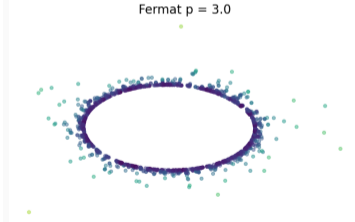
Fermat $p = 2.0$



Fermat $p = 2.5$



Fermat $p = 3.0$



Convergence of metric spaces

For $p > 1$ and $q = (p - 1)/d$,

- **Population metric space:** $(\mathcal{M}, d_{f,q})$;
- **Sample metric space:** $(\mathbb{X}_n, d_{\mathbb{X}_n,p})$.

Convergence of metric spaces

For $p > 1$ and $q = (p - 1)/d$,

- **Population metric space:** $(\mathcal{M}, d_{f,q})$;
- **Sample metric space:** $(\mathbb{X}_n, d_{\mathbb{X}_n,p})$.

Theorem (Borghini, F., Groisman, Mindlin, 2020)

There exist a constant $C(n, p, d) > 0$ such that for every $\lambda \in ((p - 1)/pd, 1/d)$ and $\varepsilon > 0$ there exist $\theta > 0$ satisfying

$$\mathbb{P} \left(d_{GH} \left((\mathcal{M}, d_{f,q}), (\mathbb{X}_n, C(n, p, d)d_{\mathbb{X}_n,p}) \right) > \varepsilon \right) \leq \exp \left(-\theta n^{(1-\lambda d)/(d+2p)} \right)$$

for n large enough.

Convergence of persistence diagrams

For $p > 1$ and $q = (p - 1)/d$,

- **Population persistence diagram:** $\text{dgm}(\text{Filt}(\mathcal{M}, d_{f,q}))$;
- **Sample persistence diagram:** $\text{dgm}(\text{Filt}(\mathbb{X}_n, d_{\mathbb{X}_n,p}))$.

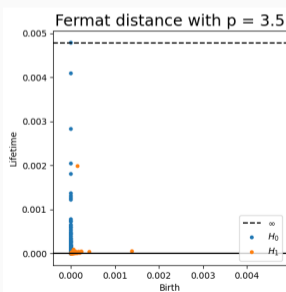
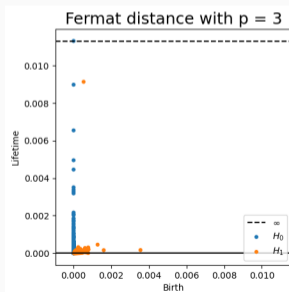
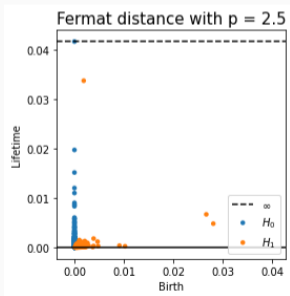
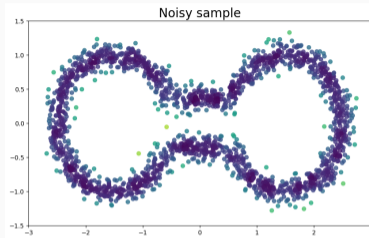
Corollary (Borghini, F., Groisman, Mindlin, 2020)

There exist a constant $C(n, p, d)$ such that for every $\lambda \in ((p - 1)/pd, 1/d)$ and $\varepsilon > 0$ there exist $\theta > 0$ satisfying

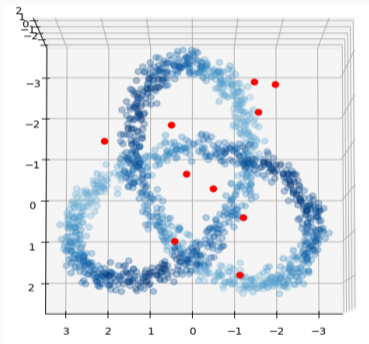
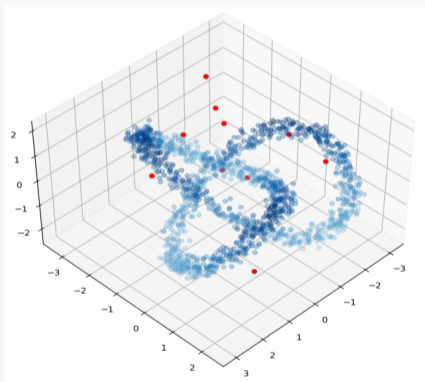
$$\begin{aligned} \mathbb{P}\left(d_b(\text{dgm}(\text{Filt}(\mathcal{M}, d_{f,q})), \text{dgm}(\text{Filt}(\mathbb{X}_n, C(n, p, d)d_{\mathbb{X}_n,p}))) > \varepsilon\right) \\ \leq \exp\left(-\theta n^{(1-\lambda d)/(d+2p)}\right) \end{aligned}$$

for n large enough.

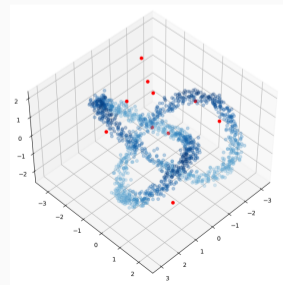
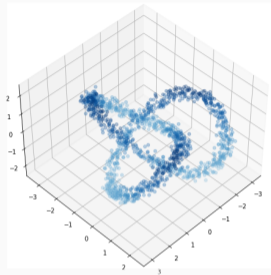
Example



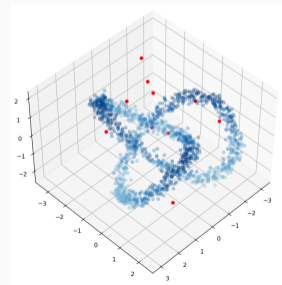
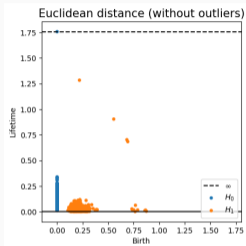
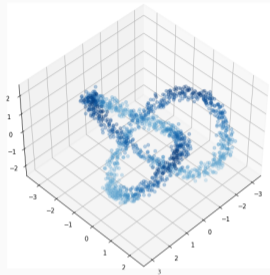
Robustness to outliers



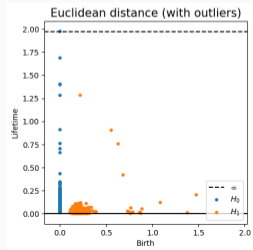
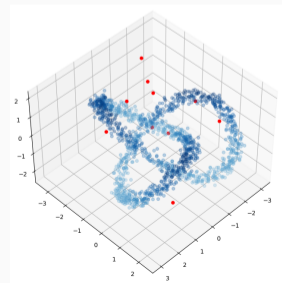
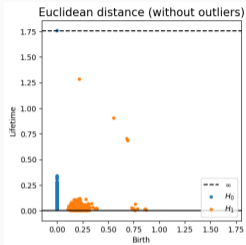
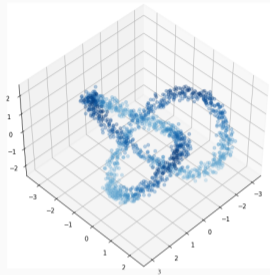
Robustness to outliers



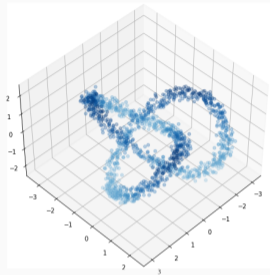
Robustness to outliers



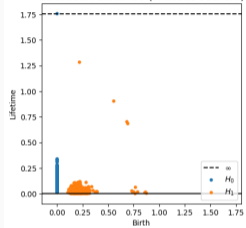
Robustness to outliers



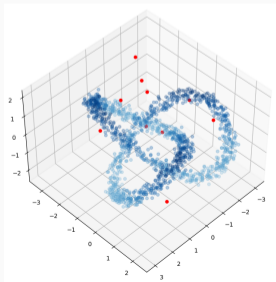
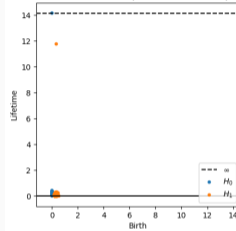
Robustness to outliers



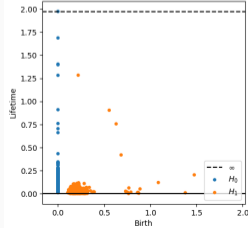
Euclidean distance (without outliers)



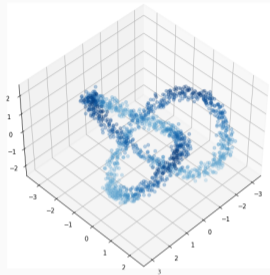
k-NN distance (without outliers)



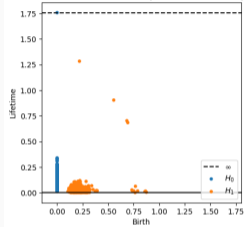
Euclidean distance (with outliers)



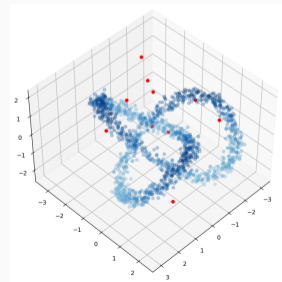
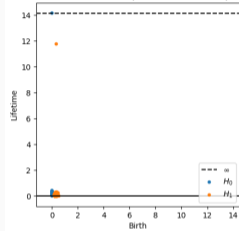
Robustness to outliers



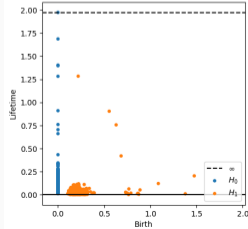
Euclidean distance (without outliers)



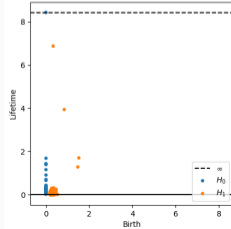
k-NN distance (without outliers)



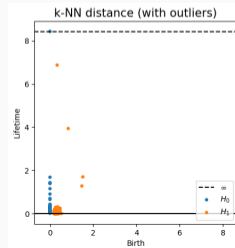
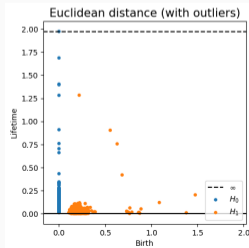
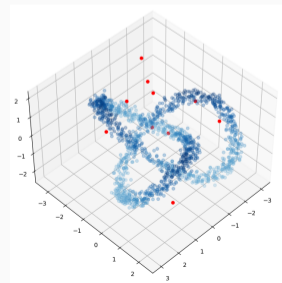
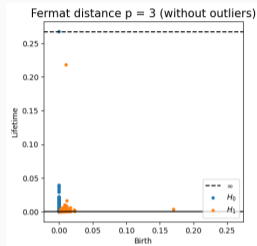
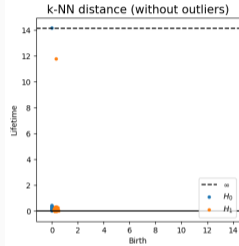
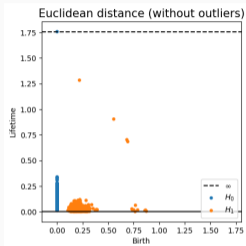
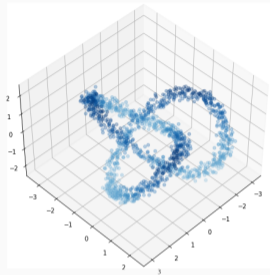
Euclidean distance (with outliers)



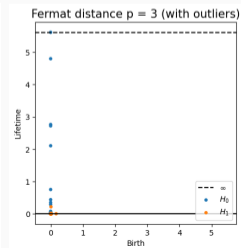
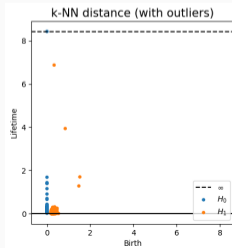
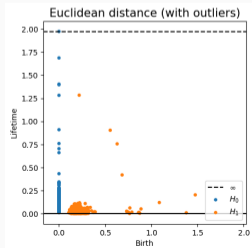
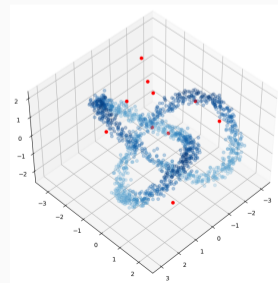
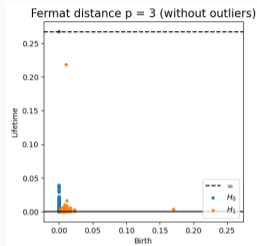
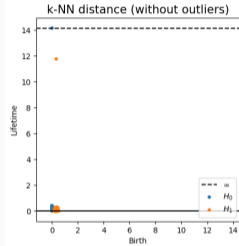
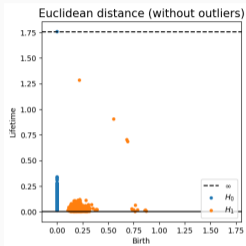
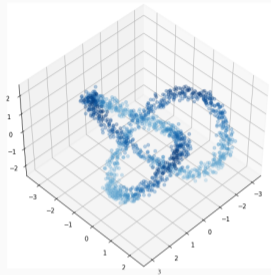
k-NN distance (with outliers)



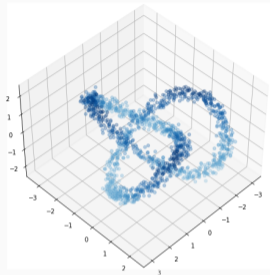
Robustness to outliers



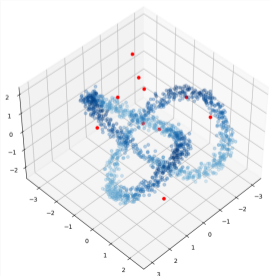
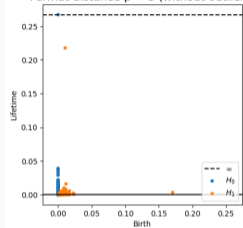
Robustness to outliers



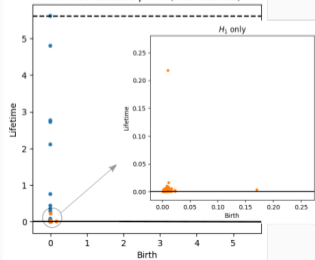
Robustness to outliers (Fermat distance)



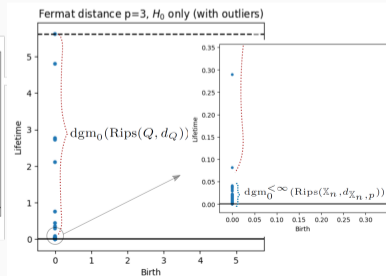
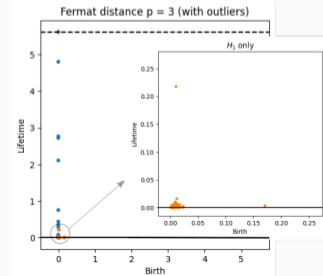
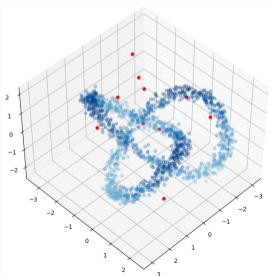
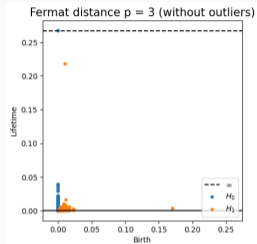
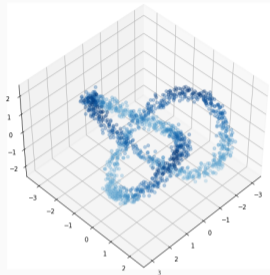
Fermat distance $p = 3$ (without outliers)



Fermat distance $p = 3$ (with outliers)



Robustness to outliers (Fermat distance)



Proposition (Borghini, F., Groisman, Mindlin, 2021)

Let \mathbb{X}_n be **sample** of \mathcal{M} and let $Y \subseteq \mathbb{R}^D \setminus \mathcal{M}$ be a finite set of **outliers**.

There exists $\delta > 0$ such that for all $k > 0$ and $p > 1$,

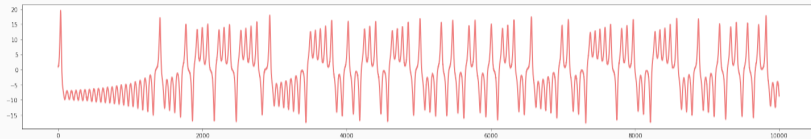
$$\text{dgm}_k(\text{Rips}_{<\delta^p}(\mathbb{X}_n \cup Y, d_{\mathbb{X}_n \cup Y, p})) = \text{dgm}_k(\text{Rips}_{<\delta^p}(\mathbb{X}_n, d_{\mathbb{X}_n, p})).$$

Here, for p large enough $\delta^p > \text{diam}(\mathbb{X}_n, d_{\mathbb{X}_n, p})$.

Applications to signal analysis

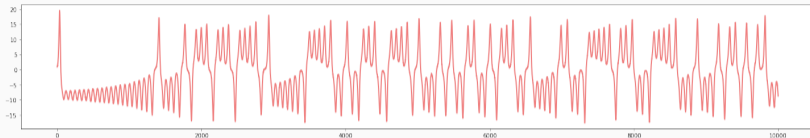
Delay embedding

- Signal $X : [t_0, t_1] \rightarrow \mathbb{R}$



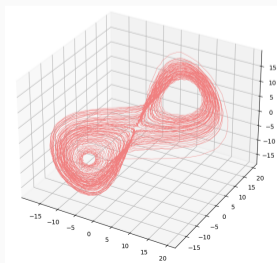
Delay embedding

- Signal $X : [t_0, t_1] \rightarrow \mathbb{R}$



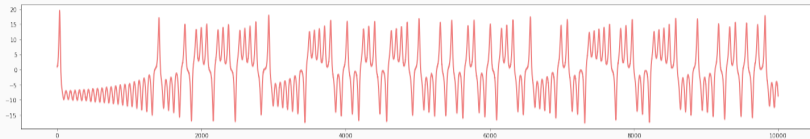
- Delay embedding

$$\mathcal{M} = \{(X(t), X(t+T), X(t+2T), \dots, X(t+(D-1)T)) : t \in [t_0, t_1 - (D-1)T]\}$$



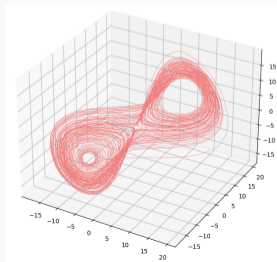
Delay embedding

- Signal $X : [t_0, t_1] \rightarrow \mathbb{R}$

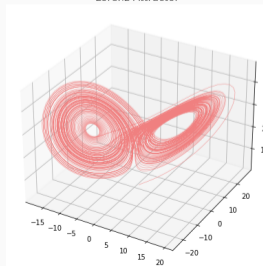


- Delay embedding

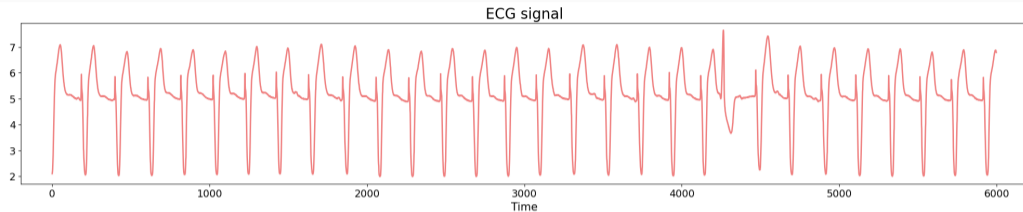
$$\mathcal{M} = \{(X(t), X(t+T), X(t+2T), \dots, X(t+(D-1)T)) : t \in [t_0, t_1 - (D-1)T]\}$$



Lorenz Attractor

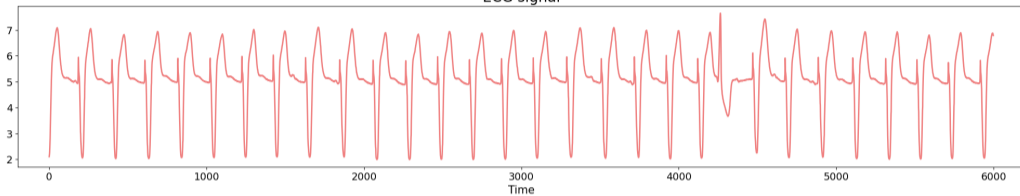


Time series: Anomaly detection

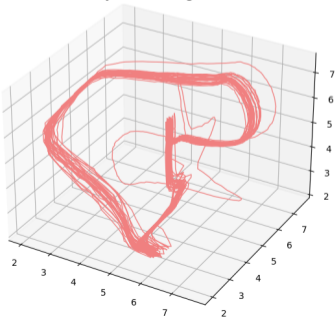


Time series: Anomaly detection

ECG signal

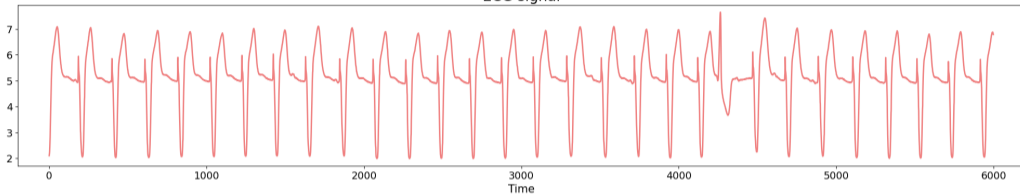


Delay Embedding $T = 15$

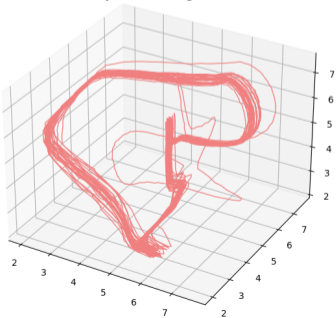


Time series: Anomaly detection

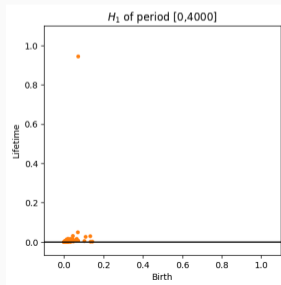
ECG signal



Delay Embedding $T = 15$

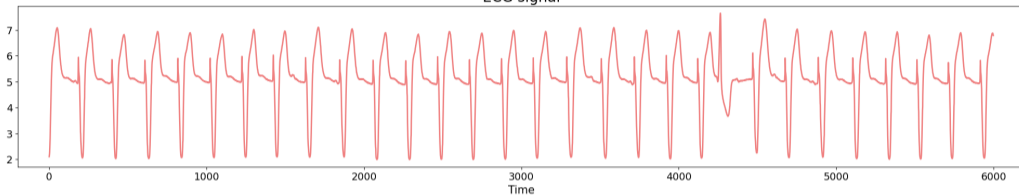


Persistence diagrams with Fermat distance for $p = 2$.

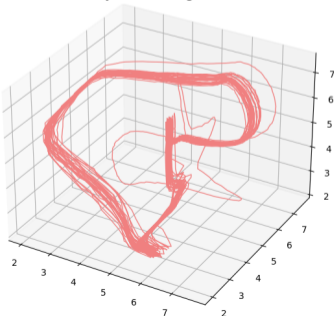


Time series: Anomaly detection

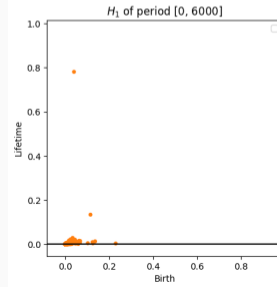
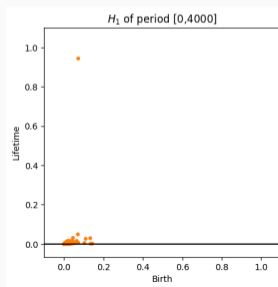
ECG signal



Delay Embedding $T = 15$

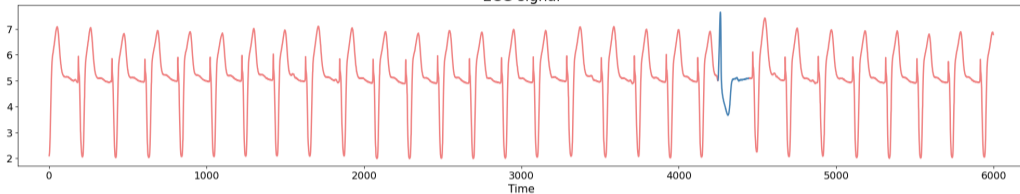


Persistence diagrams with Fermat distance for $p = 2$.

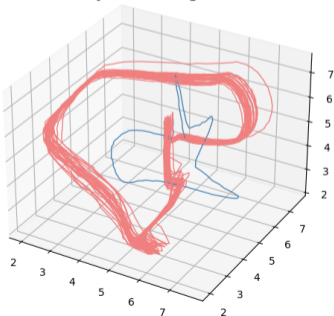


Time series: Anomaly detection

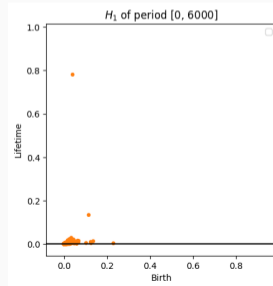
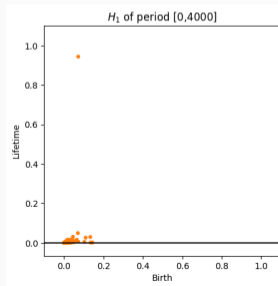
ECG signal



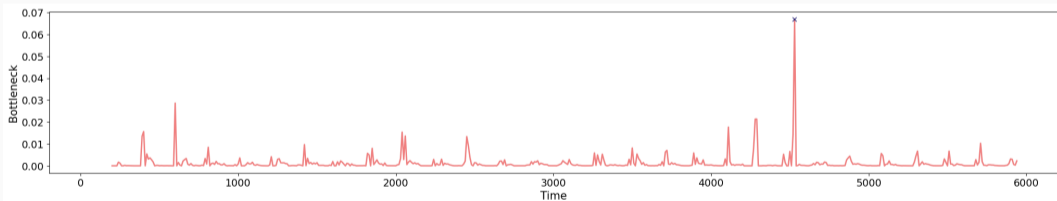
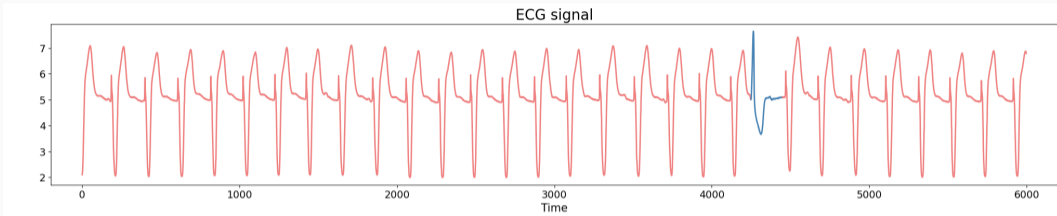
Delay Embedding $T = 15$



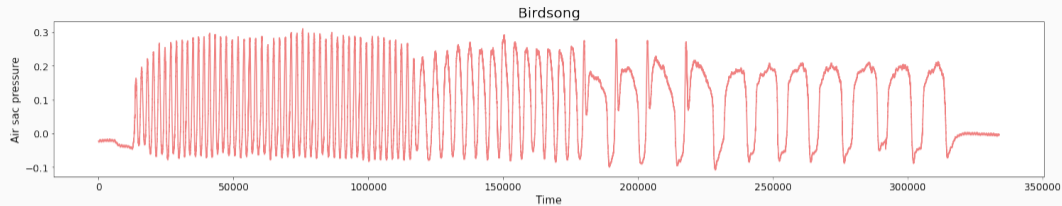
Persistence diagrams with Fermat distance for $\rho = 2$.



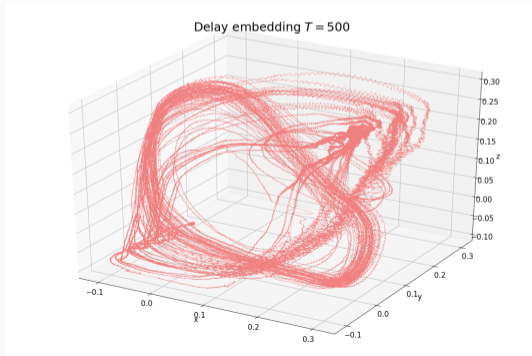
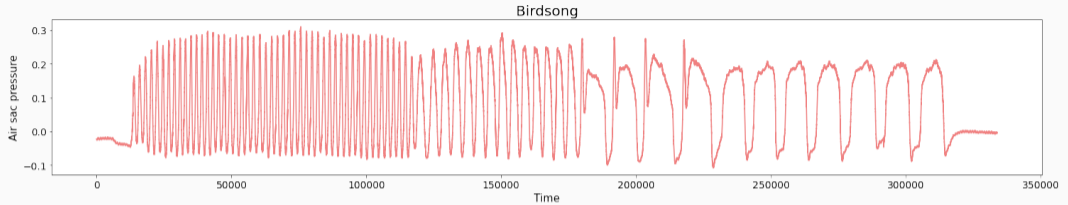
Time series: Anomaly detection



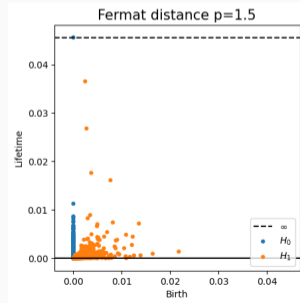
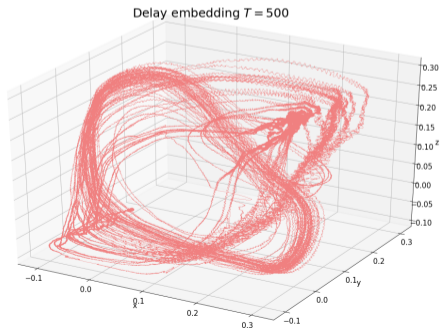
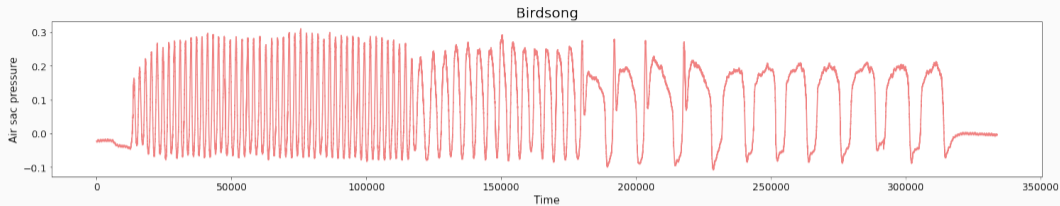
Time series: Pattern recognition



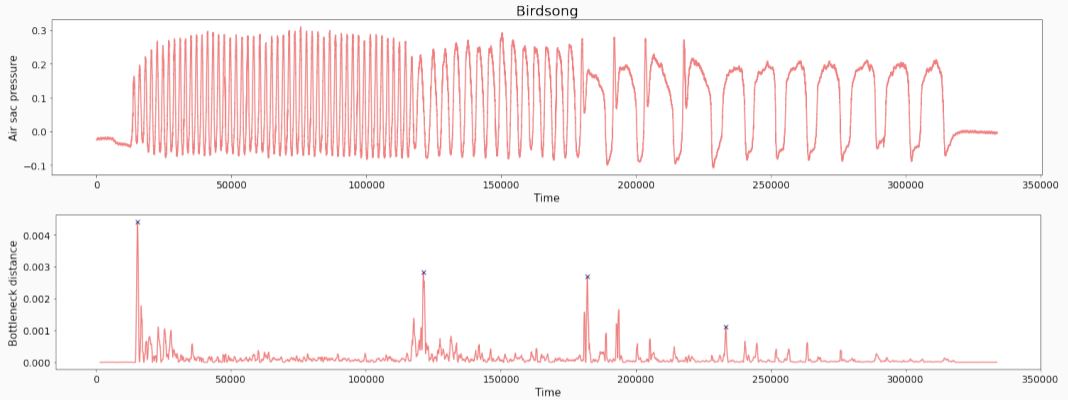
Time series: Pattern recognition



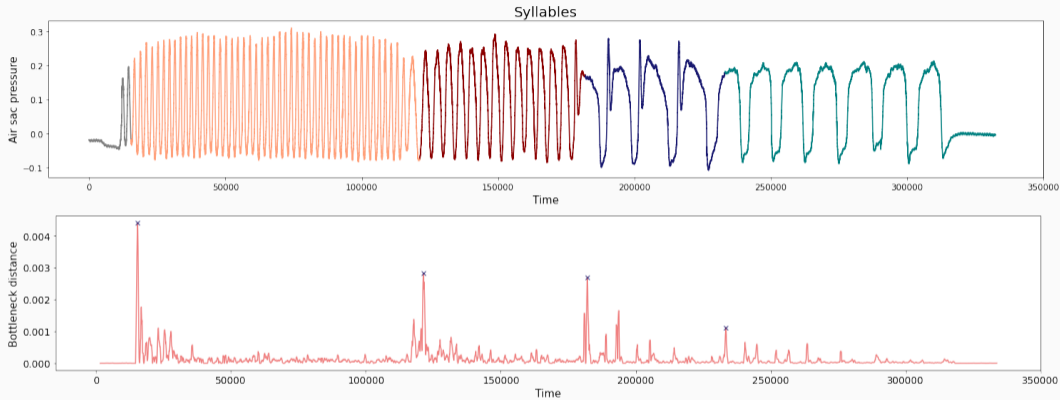
Time series: Pattern recognition



Time series: Pattern recognition

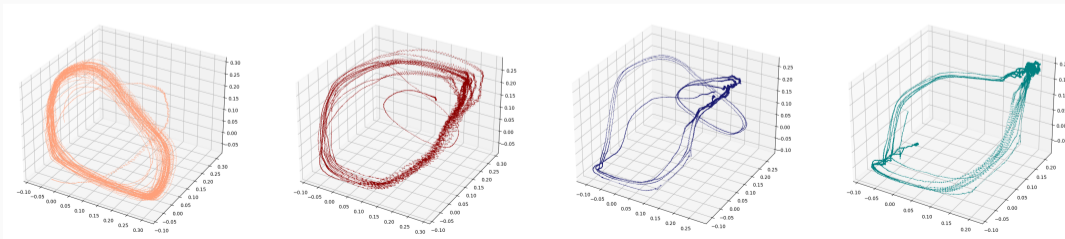
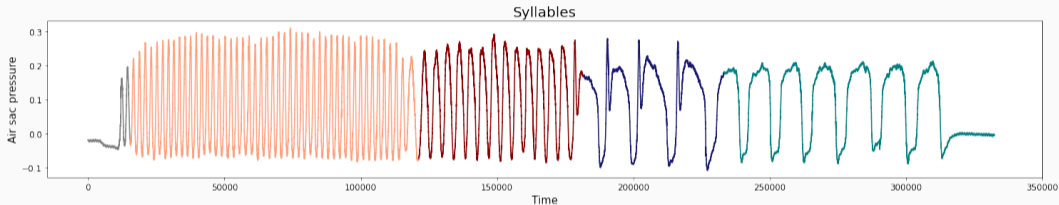


Time series: Pattern recognition



Time series: Pattern recognition

A canary song is composed by a concatenation of different syllabus patterns in the pressure in their air sacs.



- *Preprint*: E. Borghini, X. F., P. Groisman, G. Mindlin. *Intrinsic persistent homology via density-based metric learning*. arXiv:2012.07621 (2020) [Updated version soon]
- *Code*: <https://github.com/ximenafernandez/intrinsicPH>
- *Python library*: `fermat`.

email: `x.l.fernandez@swansea.ac.uk`

THANKS FOR YOUR ATTENTION!