

Workshop on Data Science
Rio de Janeiro - April 2019

A short introduction to Topological Data Analysis

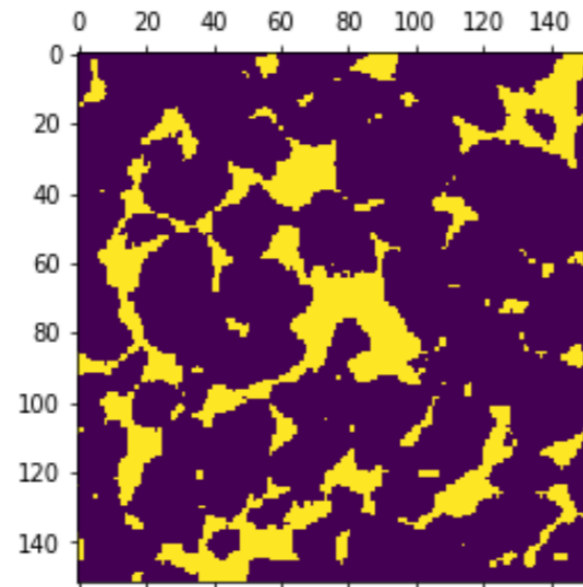
Frédéric Chazal
DataShape team
INRIA Saclay - Ile-de-France
frederic.chazal@inria.fr marc.glisse@inria.fr



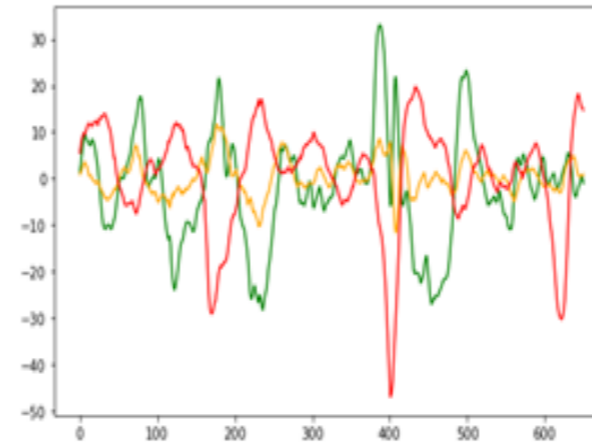
What is Topological Data Analysis (TDA)?



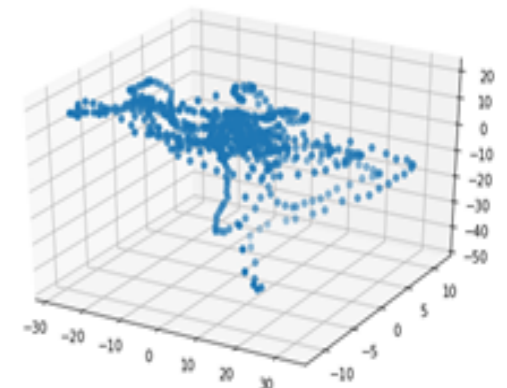
[Scanned 3D object]



[3D images (porous rocks)]

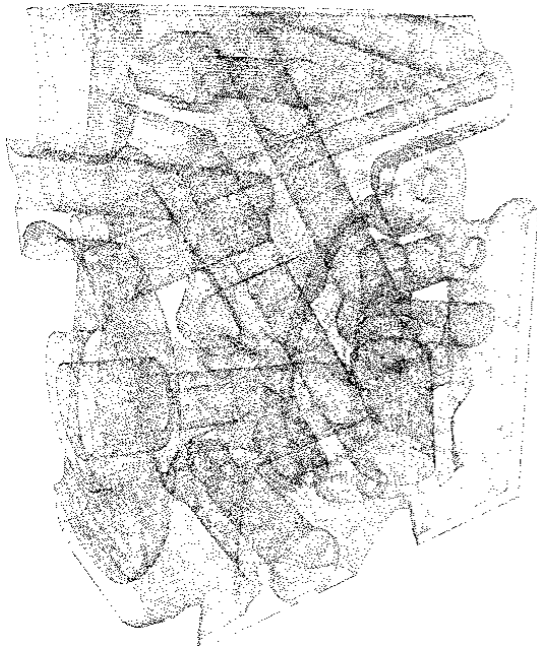


[Sensors (Sysnav courtesy)]

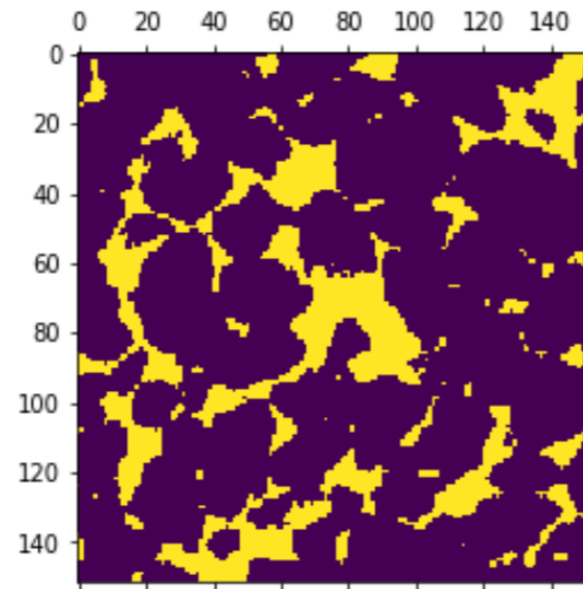


Modern data carry complex, but important, geometric/topological structure!

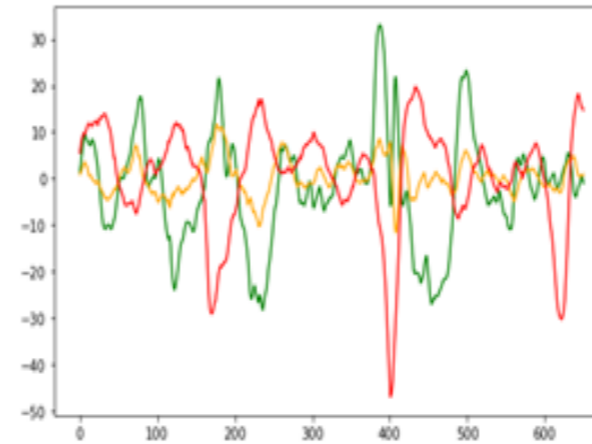
What is Topological Data Analysis (TDA)?



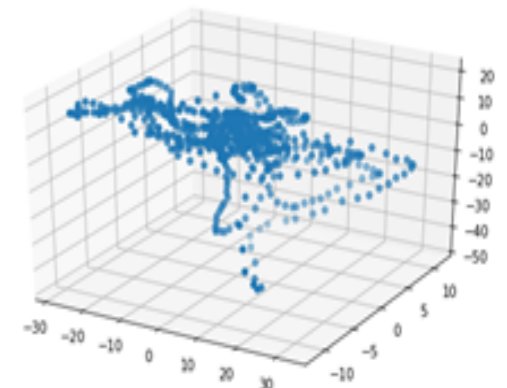
[Scanned 3D object]



[3D images (porous rocks)]



[Sensors (Sysnav courtesy)]



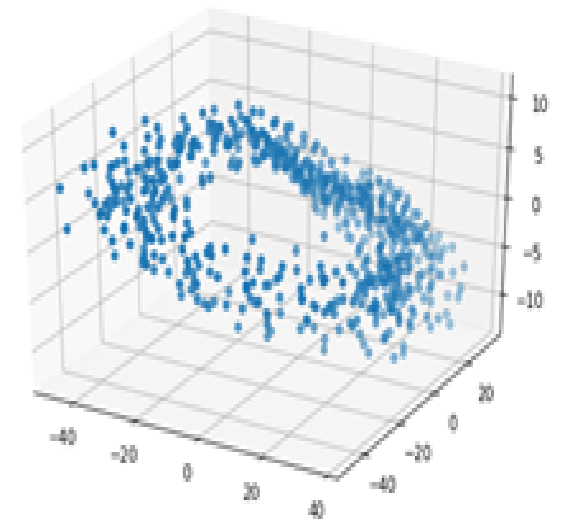
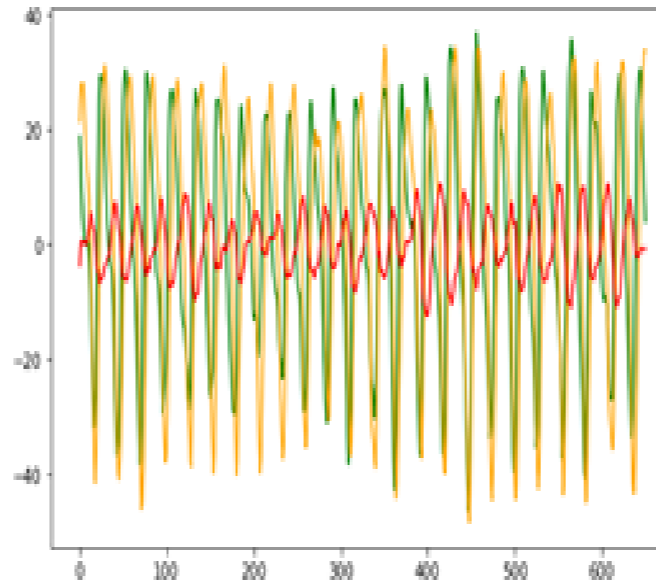
Topological Data Analysis (TDA) is a recent field whose aim is to:

- infer relevant topological and geometric features from complex data,
- take advantage of topological/geometric information for further Data Analysis, Machine Learning and AI tasks.

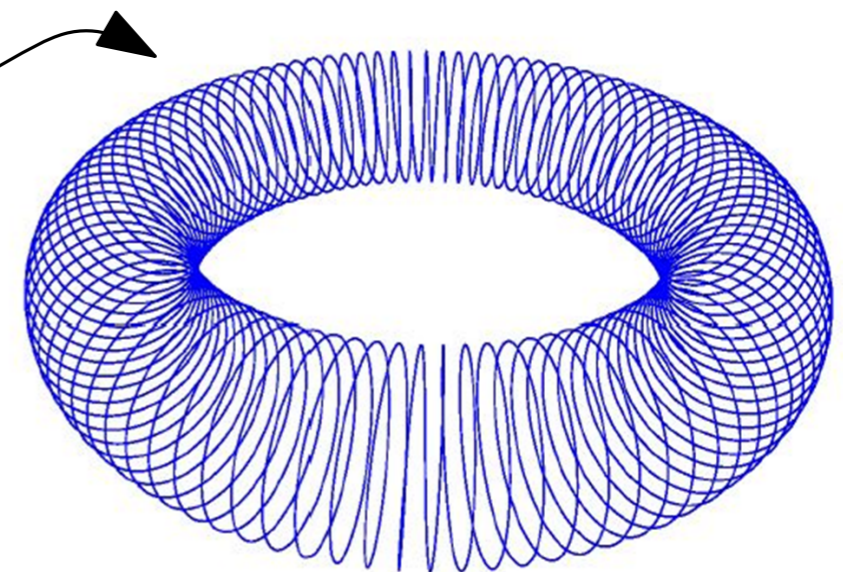
Challenges and goals

Problem(s):

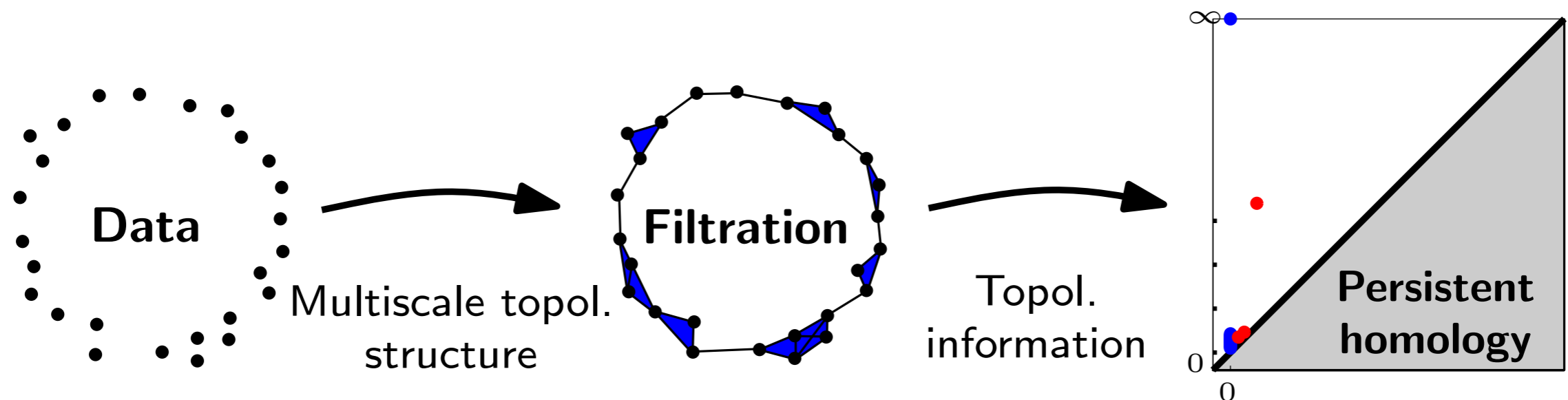
- what is topological structure of data?
- how to compare topological properties (invariants) of close shapes/data sets?



- Challenges and goals:
 - no direct access to topological/geometric information: need of intermediate constructions (simplicial complexes);
 - distinguish topological “signal” from noise;
 - topological information may be multiscale;
 - statistical analysis of topological information.

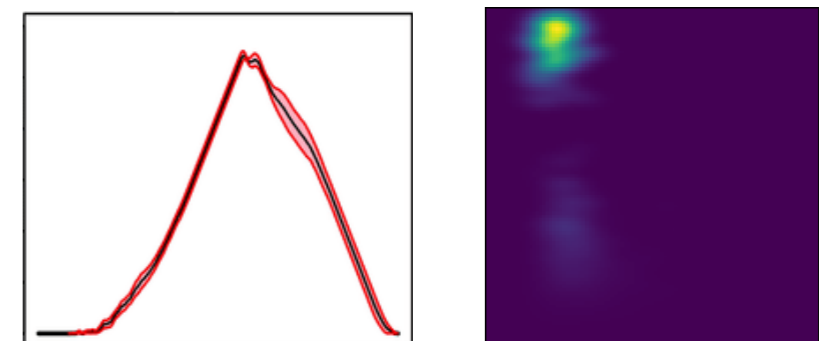


The classical TDA pipeline



1. Build a multiscale topol. structure on top of data: **filtrations**.
2. Compute multiscale topol. signatures: **persistent homology**
3. Take advantage of the signature for further Machine Learning and AI tasks: **Statistical aspects and representations of persistence**

Machine Learning / AI



Representations of persistence

Persistent homology

The theory of persistence

A recent theory that is subject to intense research activities:

- **from the mathematical perspective:**

- general algebraic framework (persistence modules) and general stability results.
- extensions and generalizations of persistence (zig-zag persistence, multi-persistence, etc...)
- Statistical analysis of persistence.

- **from the algorithmic and computational perspective:**

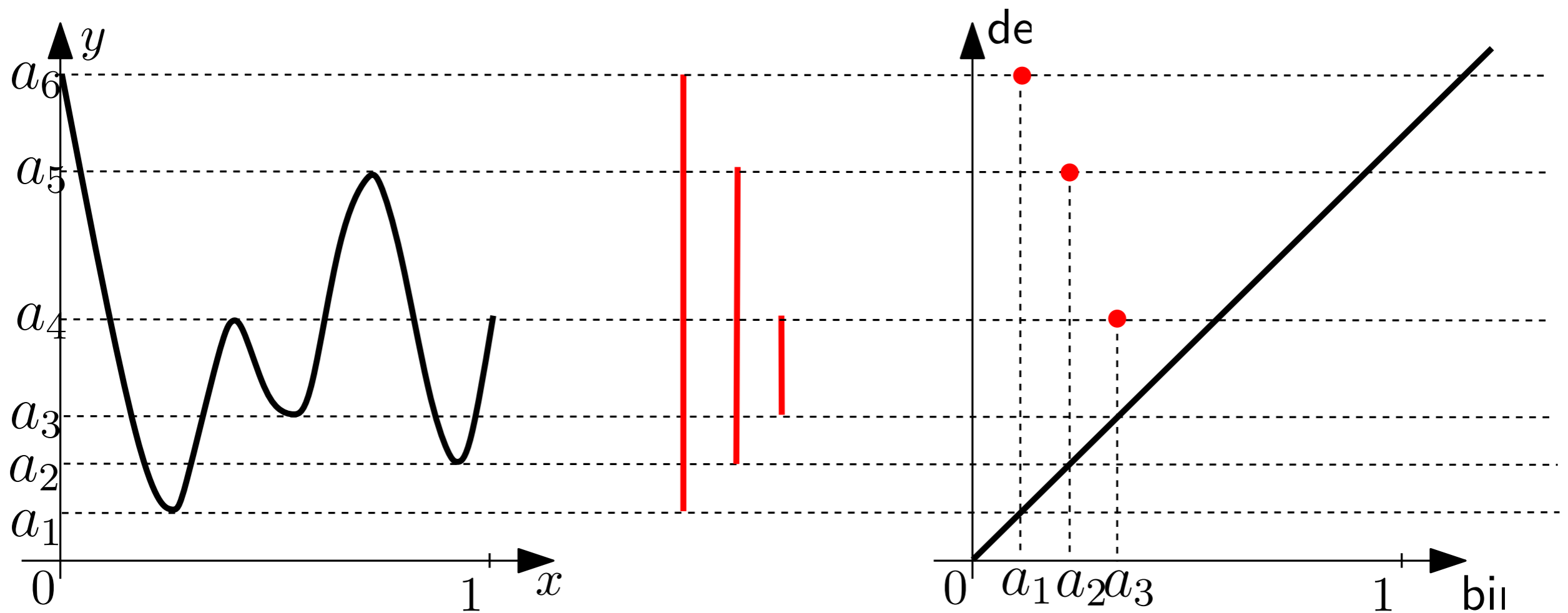
- efficient algorithms to compute persistence and some of its variants.
- efficient software libraries (in particular, Gudhi: <https://project.inria.fr/gudhi/>).

- **from the data science perspective:**

- representations of persistence that are suitable for Machine Learning
- Topological/geometric information in combination with other features

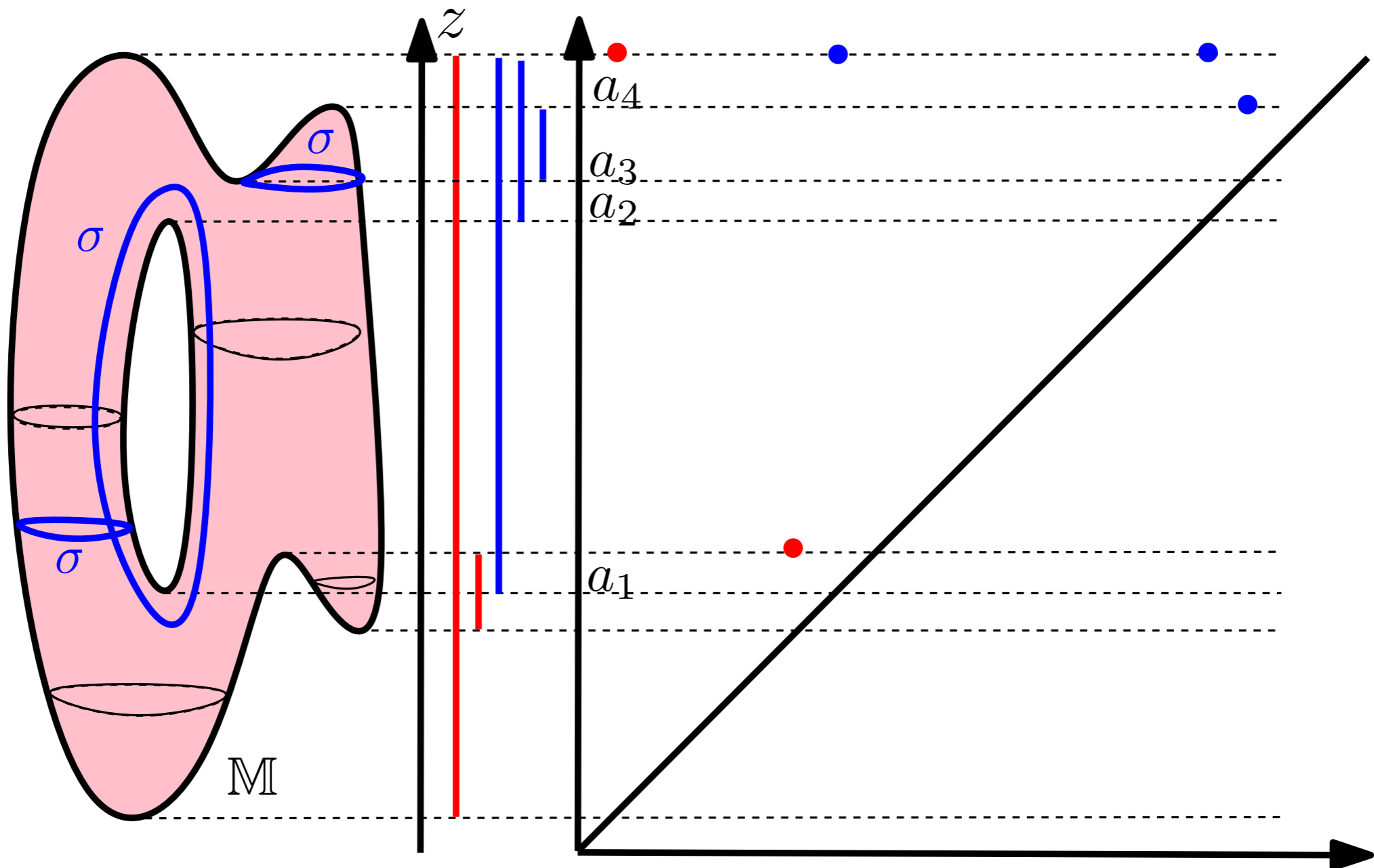
A whole machinery at the crossing of mathematics and computer science!

Persistent homology for functions



Tracking and encoding the evolution of the connected components (0-dimensional homology) of the sublevel sets of a function

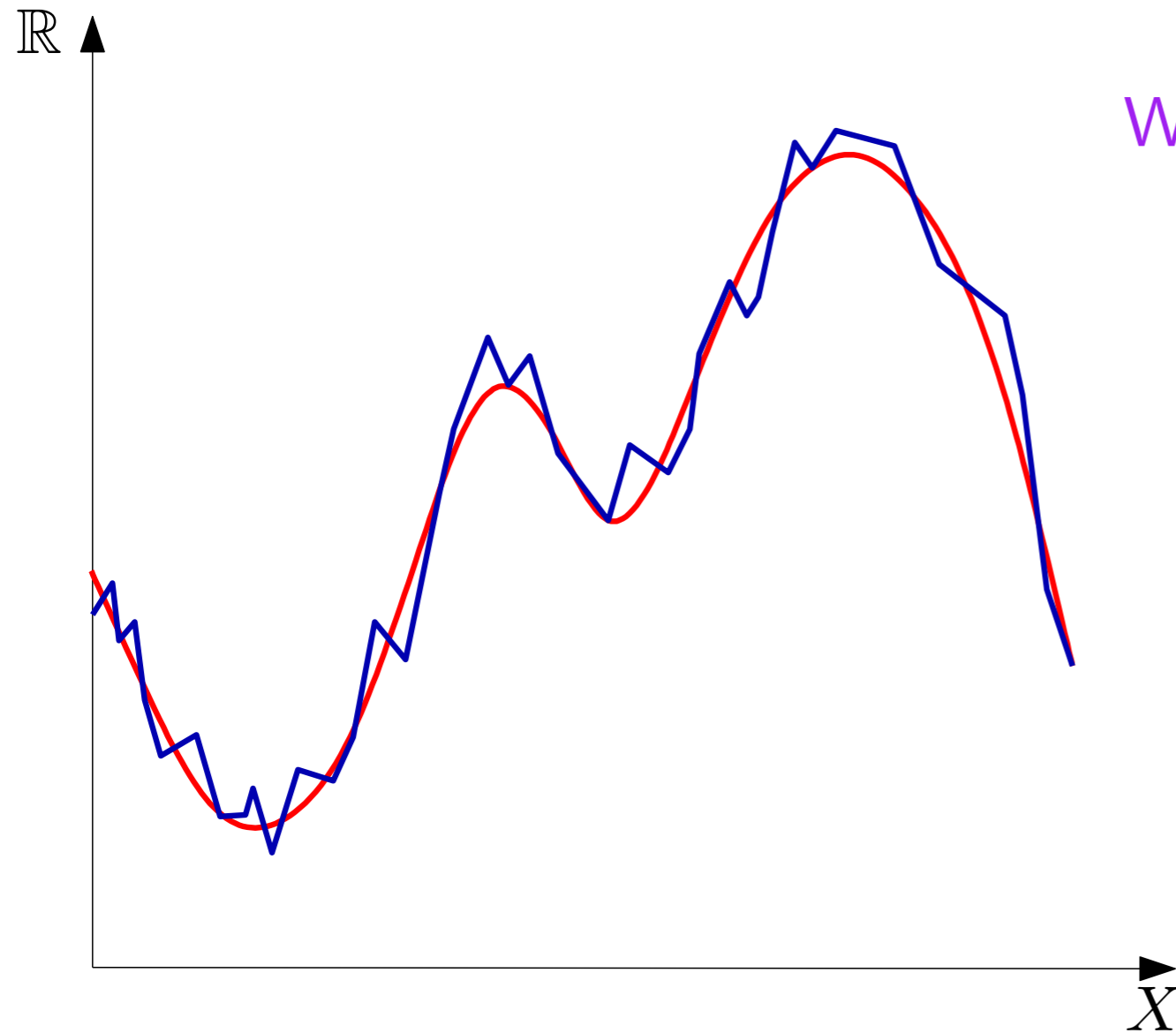
Persistent homology for functions



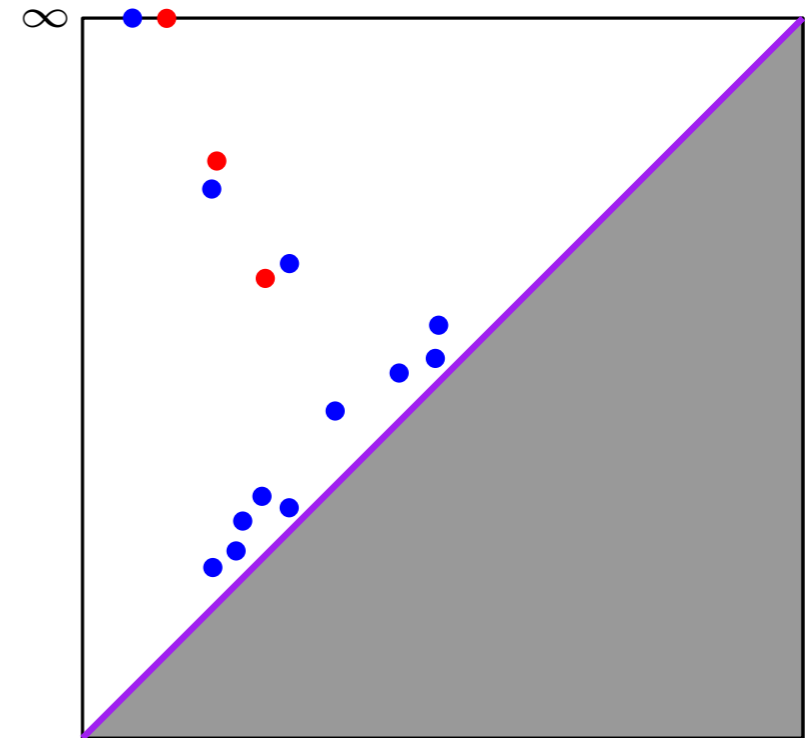
Tracking and encoding the evolution of the **connected components (0-dimensional homology)** and **cycles (1-dimensional homology)** of the sublevel sets.

Homology: an algebraic way to rigorously formalize the notion of k -dimensional cycles through a vector space (or a group), the homology group whose dimension is the number of "independent" cycles (the Betti number).

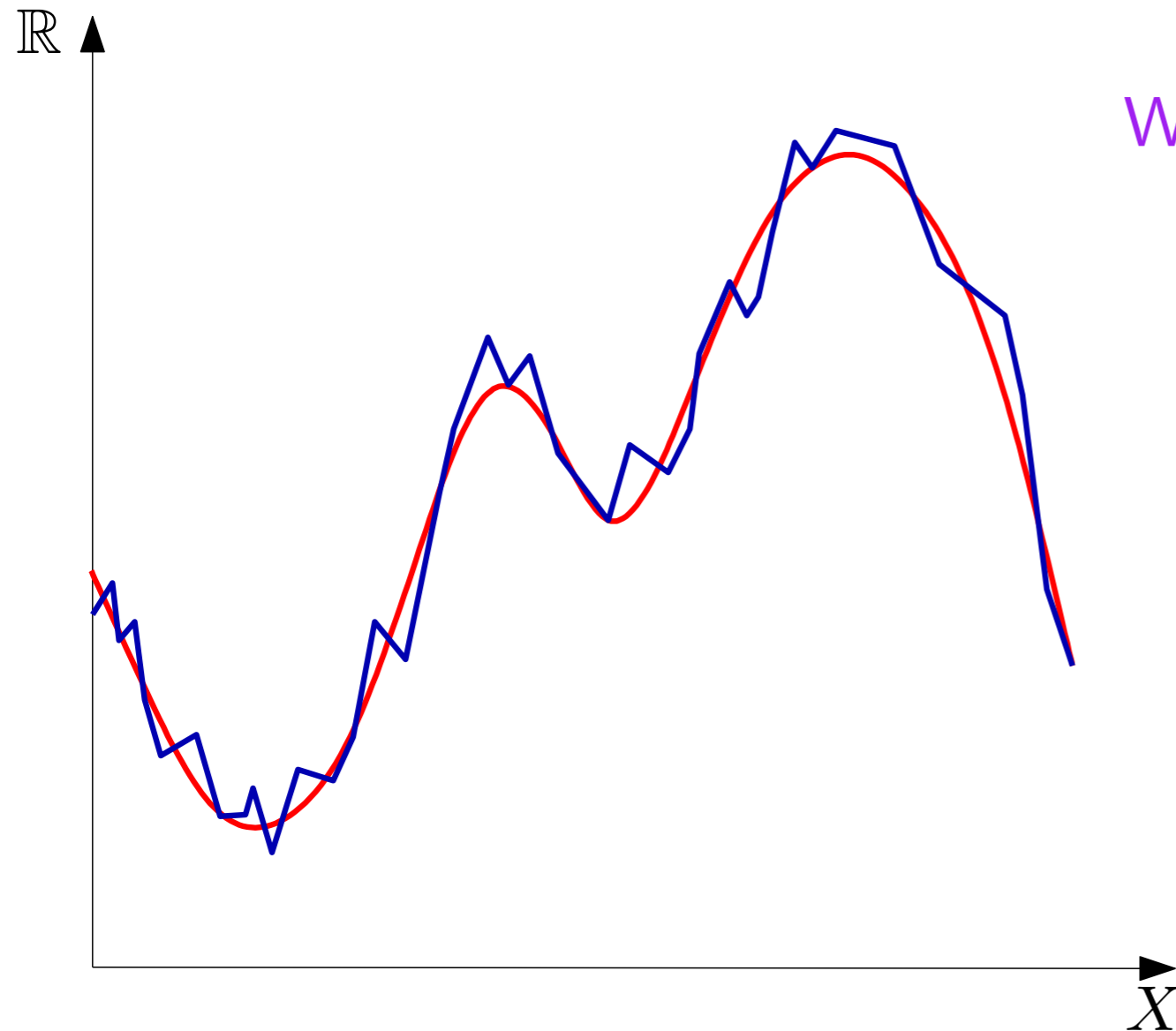
Stability properties



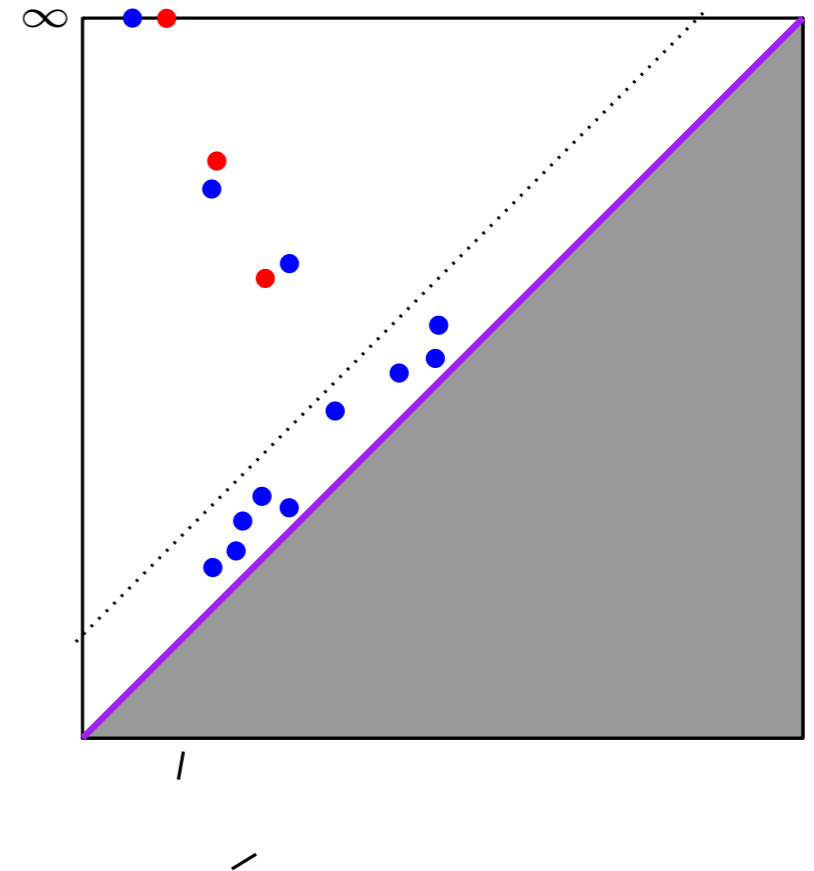
What if f is slightly perturbed?



Stability properties



What if f is slightly perturbed?

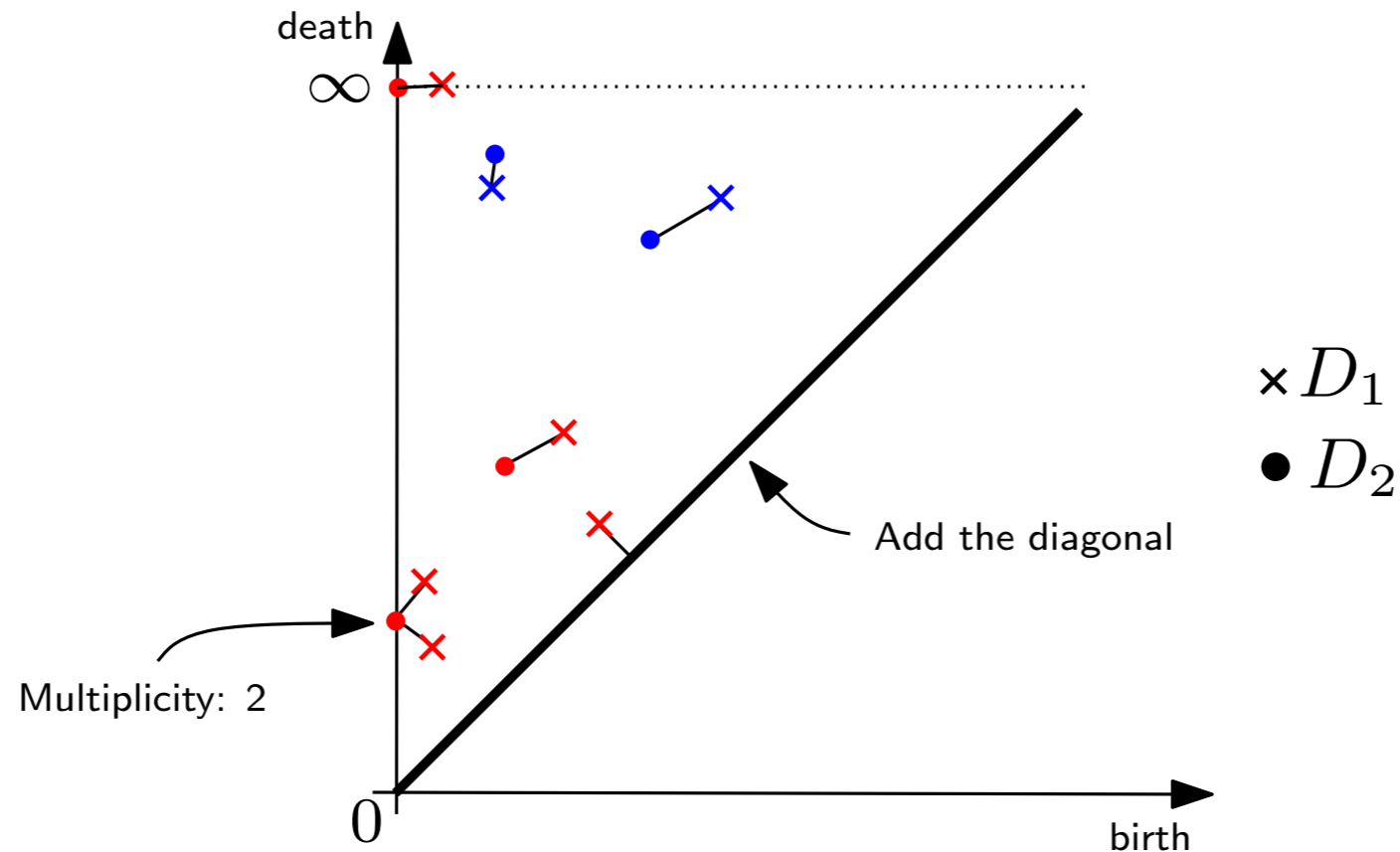


Theorem (Stability):

For any *tame* functions $f, g : \mathbb{X} \rightarrow \mathbb{R}$, $d_B(D_f, D_g) \leq \|f - g\|_\infty$.

[Cohen-Steiner, Edelsbrunner, Harer 05], [C., Cohen-Steiner, Glisse, Guibas, Oudot - SoCG 09], [C., de Silva, Glisse, Oudot 12]

Comparing persistence diagrams



The **bottleneck distance** between two diagrams D_1 and D_2 is

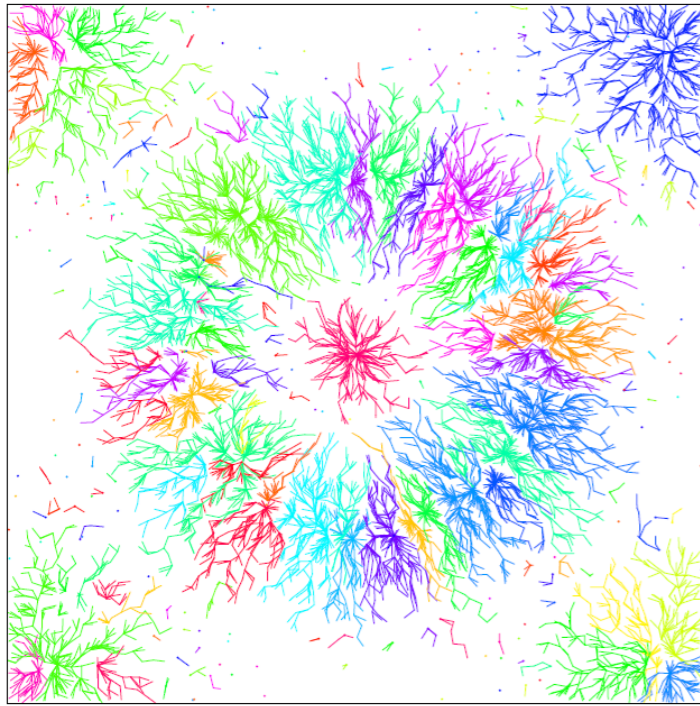
$$d_B(D_1, D_2) = \inf_{\gamma \in \Gamma} \sup_{p \in D_1} \|p - \gamma(p)\|_{\infty}$$

where Γ is the set of all the bijections between D_1 and D_2 and $\|p - q\|_{\infty} = \max(|x_p - x_q|, |y_p - y_q|)$.

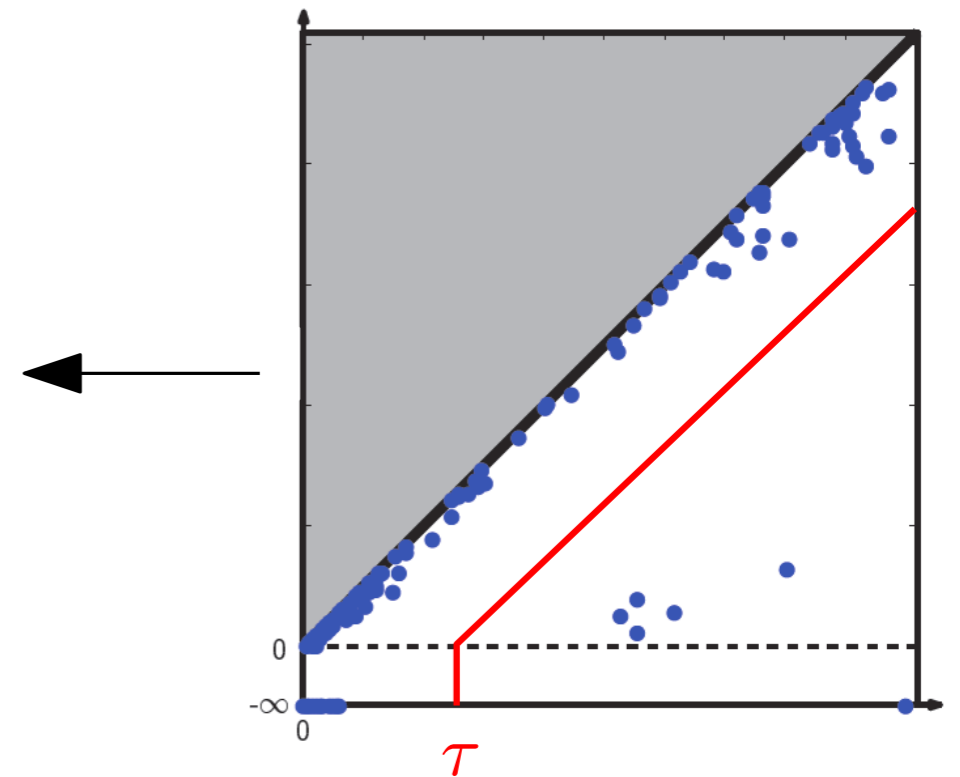
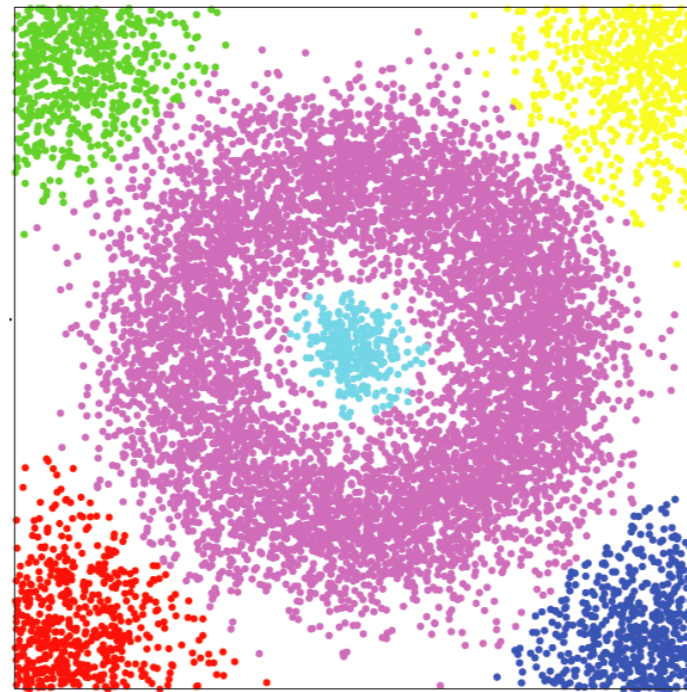
→ Persistence diagrams provide easy to compare topological signatures.

Some applications (illustrations)

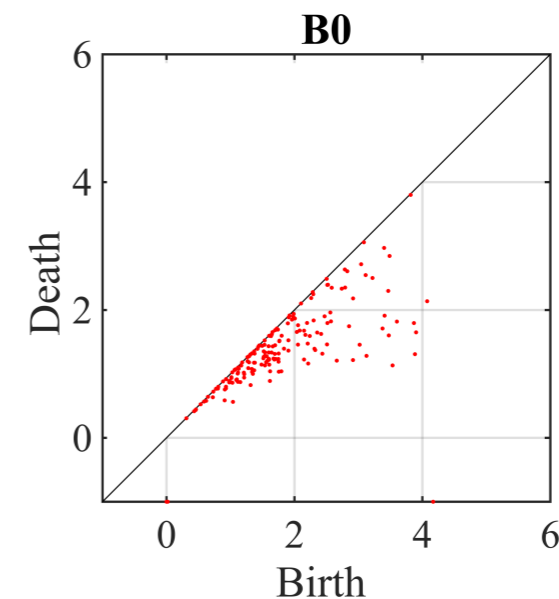
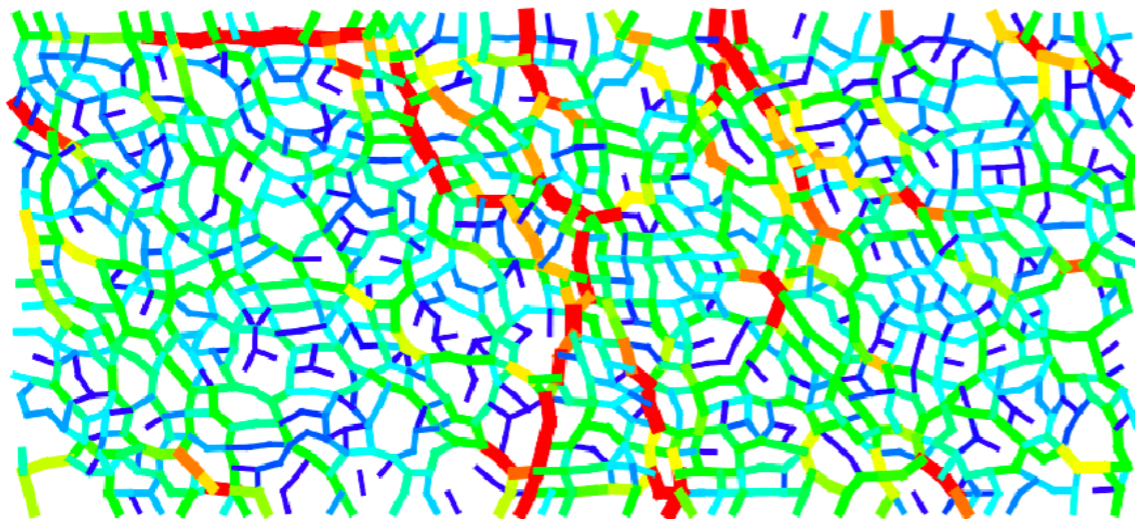
- Persistence-based clustering [C., Guibas, Oudot, Skraba - J. ACM 2013]



$\tau = 0$

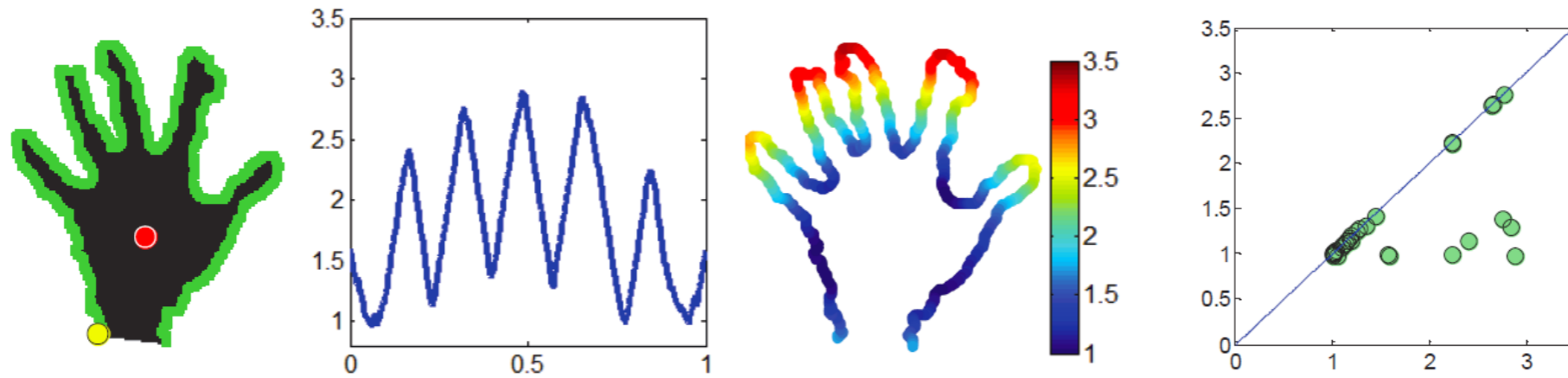


- Analysis of force fields in granular media [Kramar, Mischaikow et al]

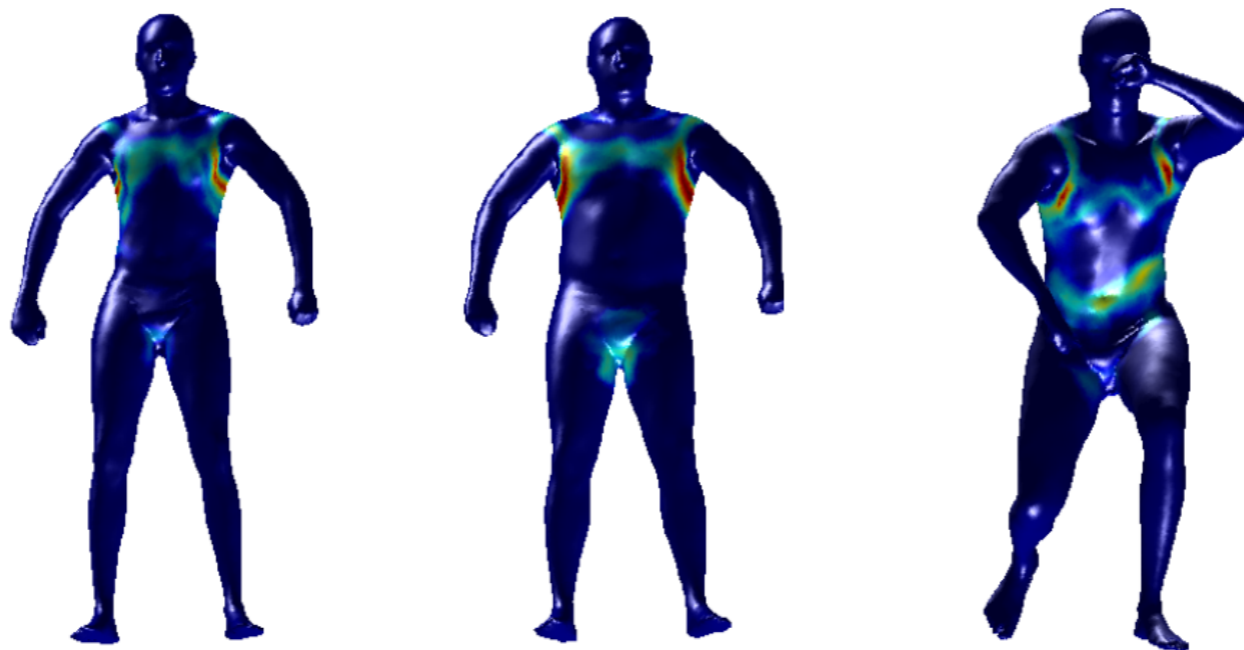


Some applications (illustrations)

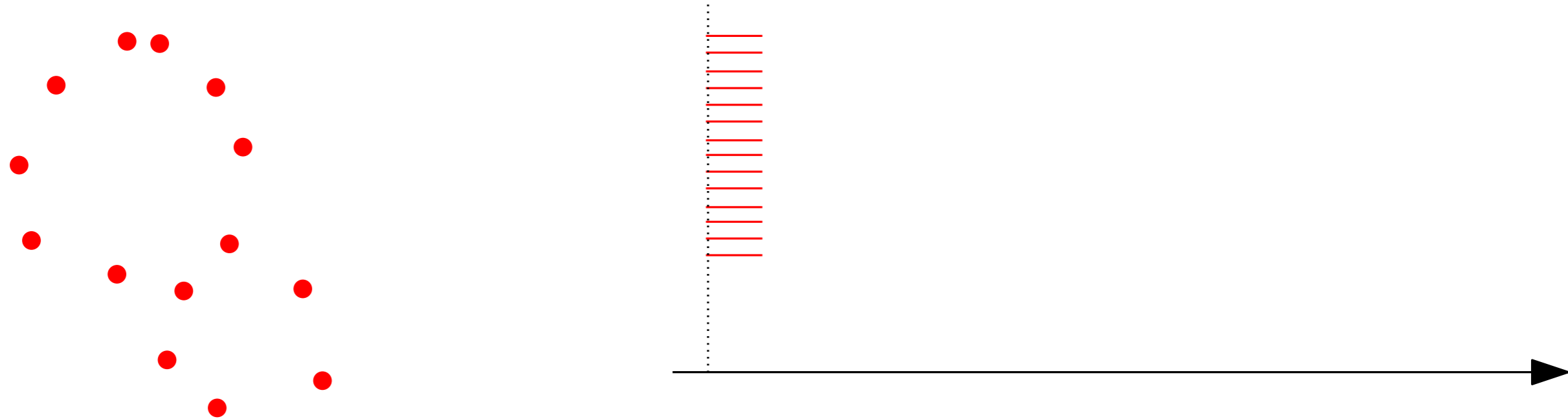
- Hand gesture recognition [Li, Ovsjanikov, C. - CVPR'14]



- Persistence-based pooling for shape recognition [Bonis, Ovsjanikov, Oudot, C. 2016]

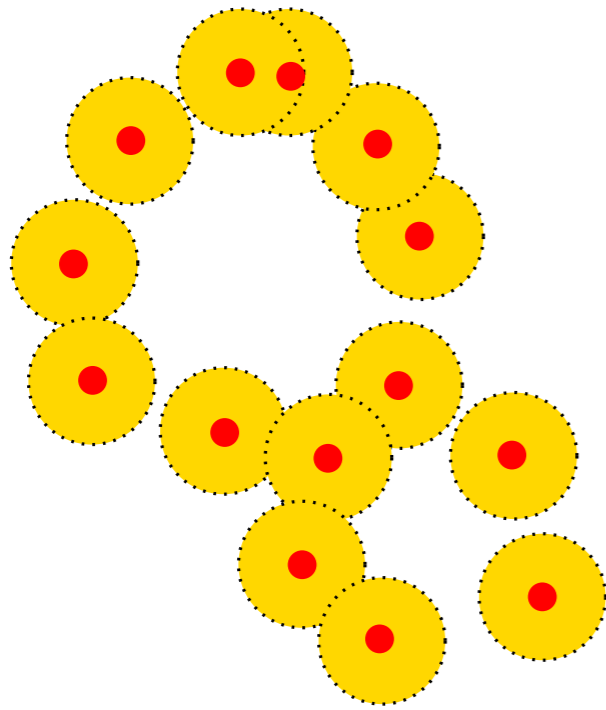


Persistent homology for point cloud data



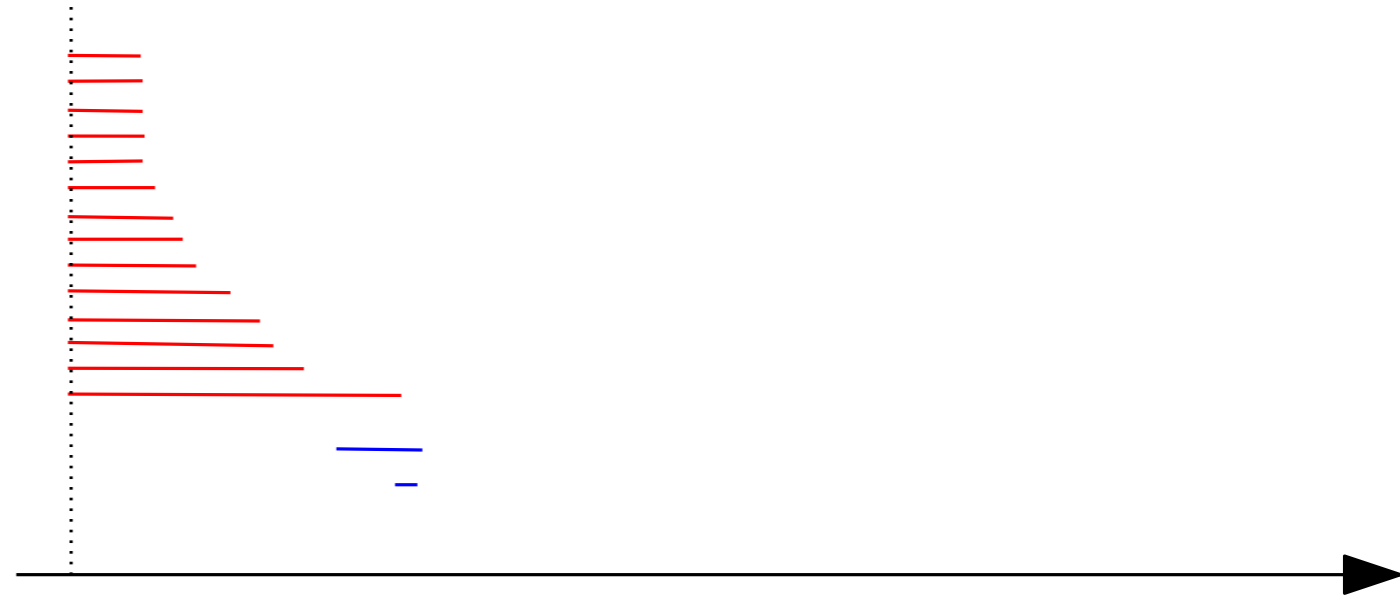
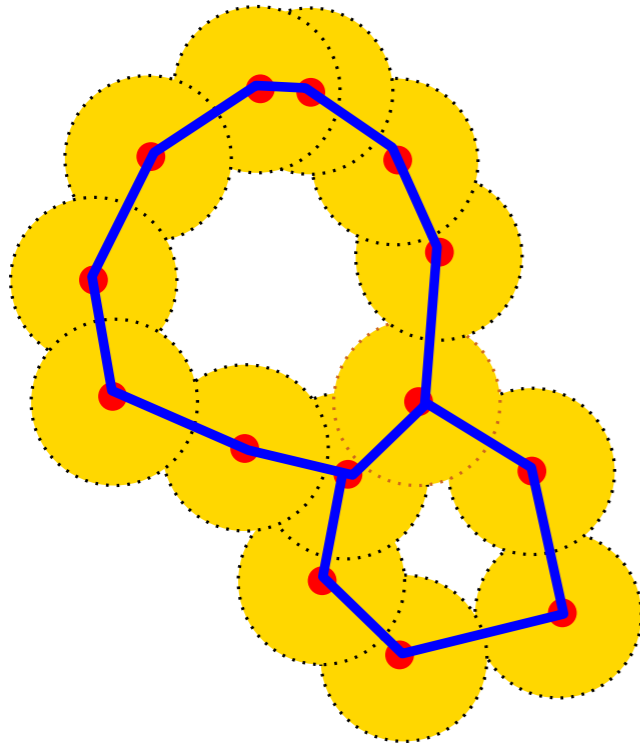
- Filtrations allow to construct “shapes” representing the data in a multiscale way.
- **Persistent homology:** encode the evolution of the topology across the scales → multi-scale topological signatures.

Persistent homology for point cloud data



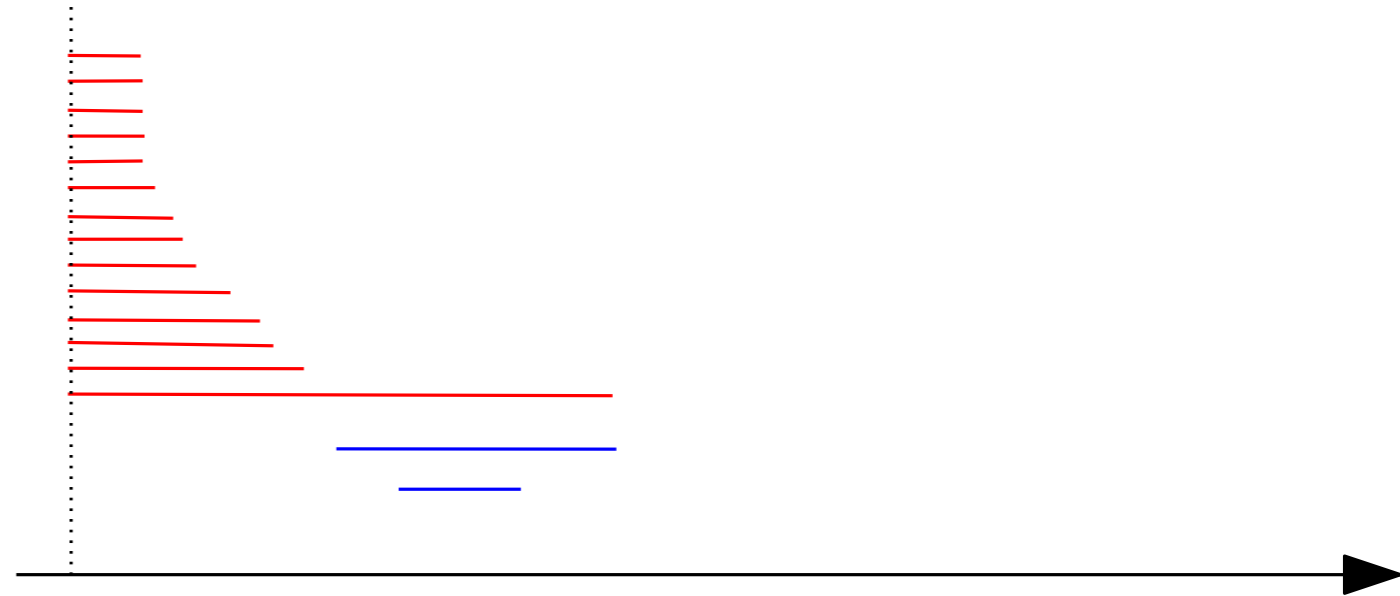
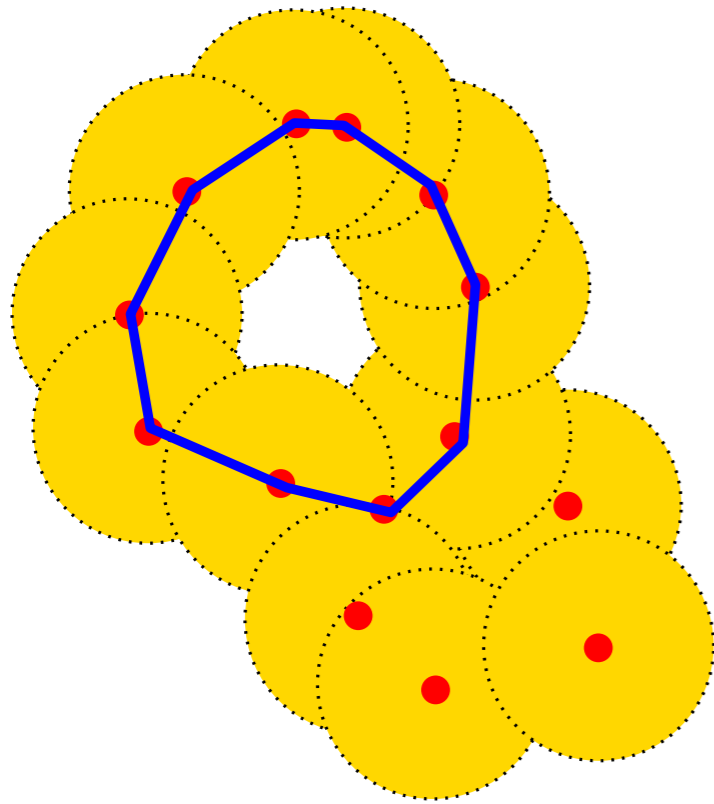
- Filtrations allow to construct “shapes” representing the data in a multiscale way.
- **Persistent homology:** encode the evolution of the topology across the scales → multi-scale topological signatures.

Persistent homology for point cloud data



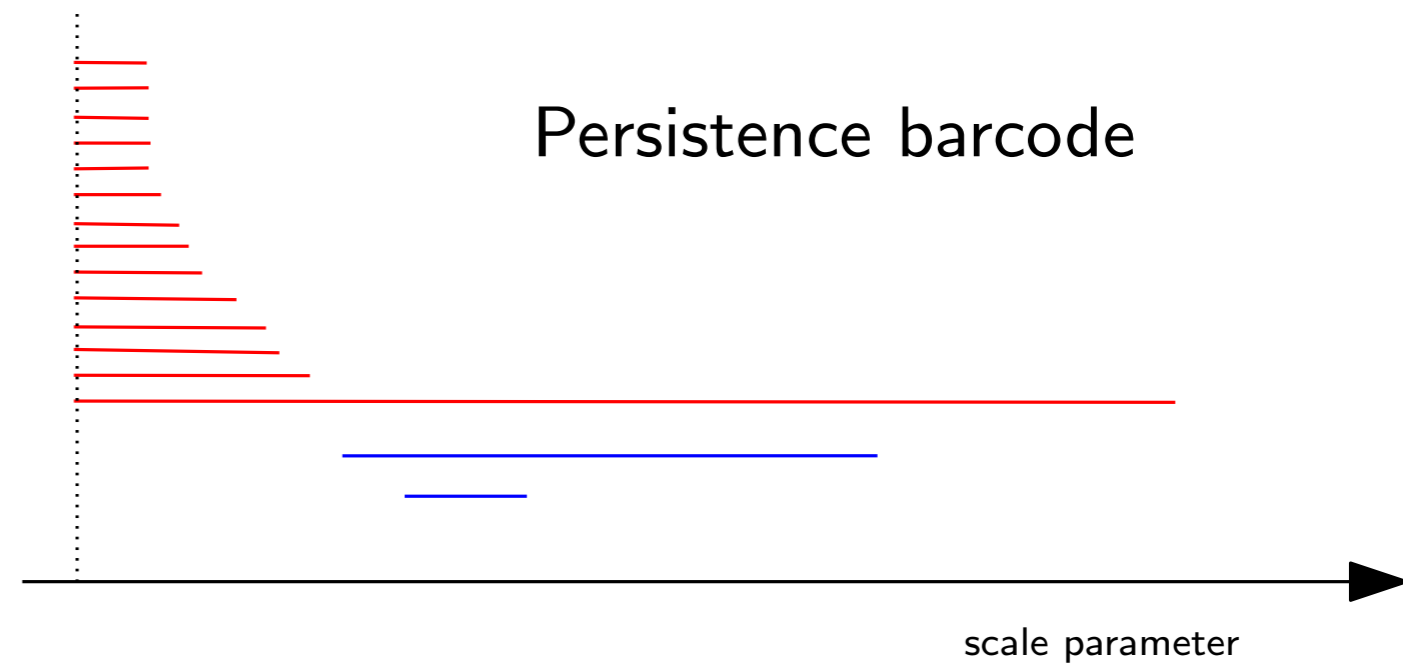
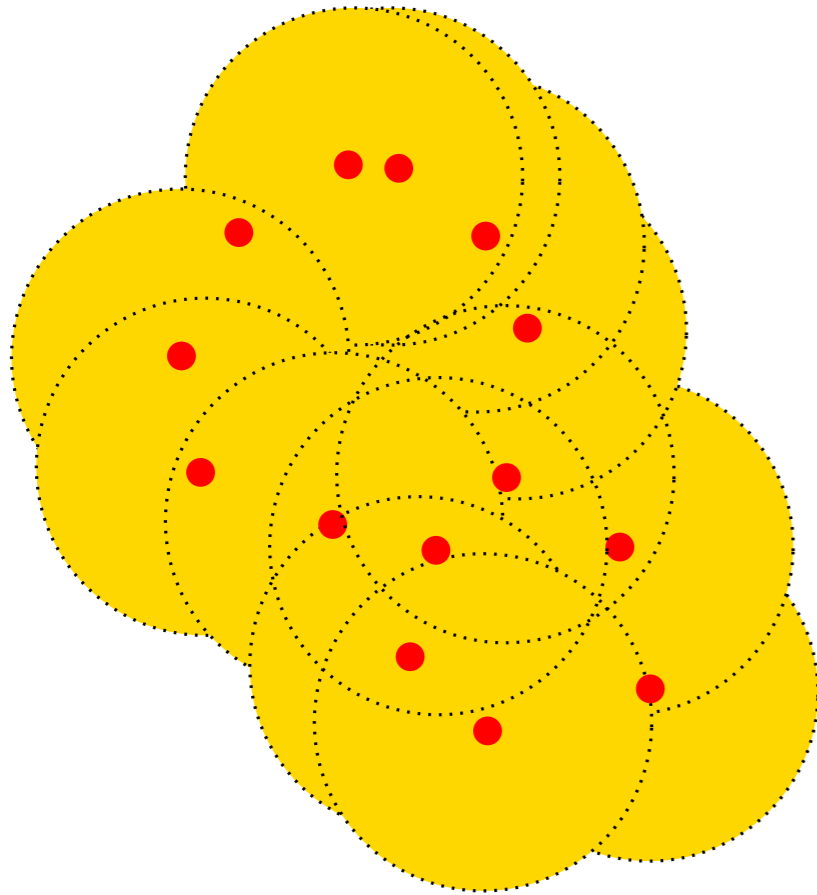
- Filtrations allow to construct “shapes” representing the data in a multiscale way.
- **Persistent homology:** encode the evolution of the topology across the scales → multi-scale topological signatures.

Persistent homology for point cloud data

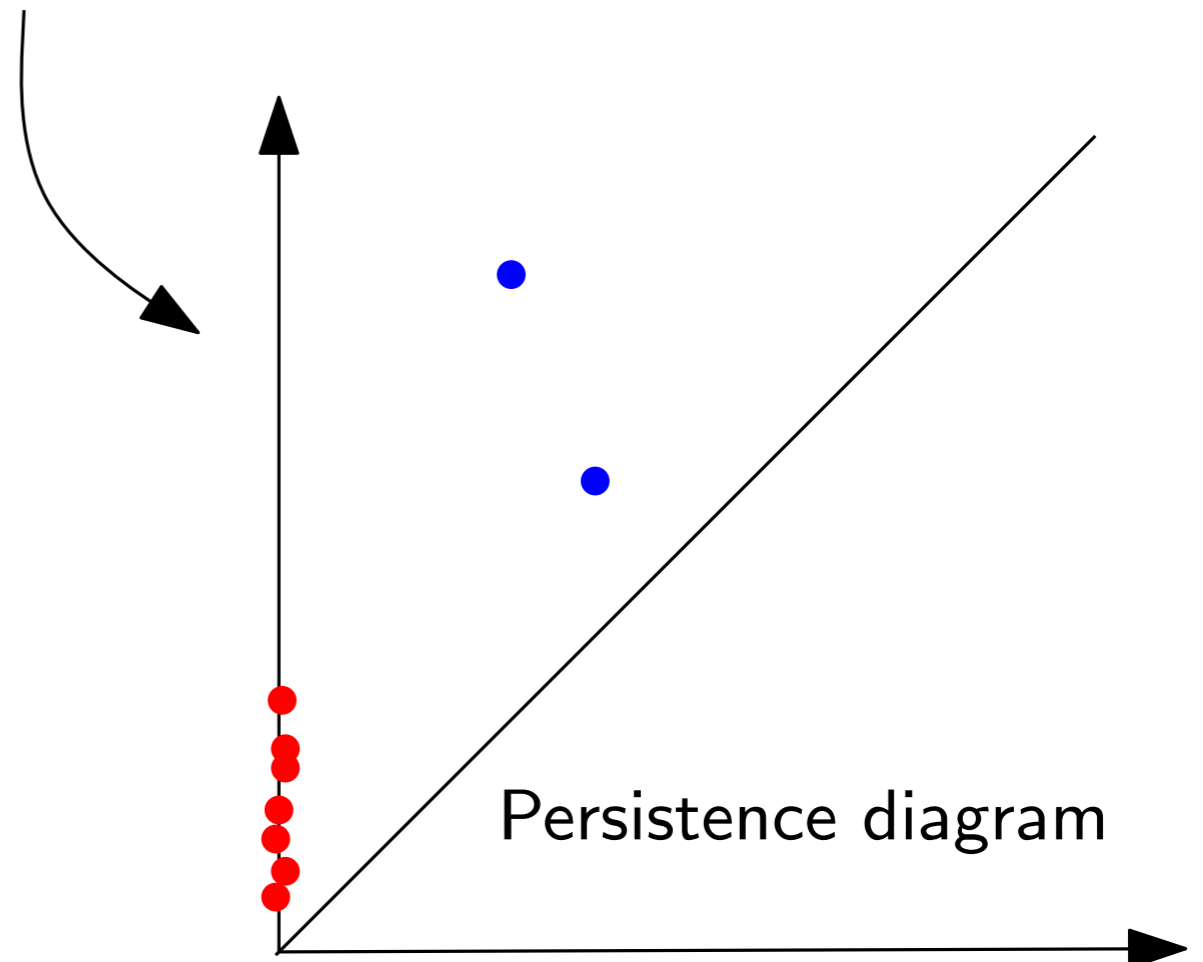


- Filtrations allow to construct “shapes” representing the data in a multiscale way.
- **Persistent homology:** encode the evolution of the topology across the scales → multi-scale topological signatures.

Persistent homology for point cloud data



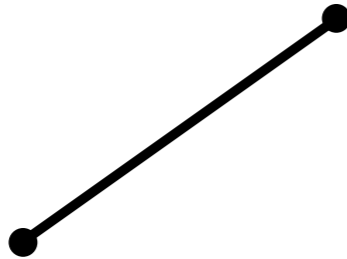
- Filtrations allow to construct “shapes” representing the data in a multiscale way.
- **Persistent homology:** encode the evolution of the topology across the scales → multi-scale topological signatures.



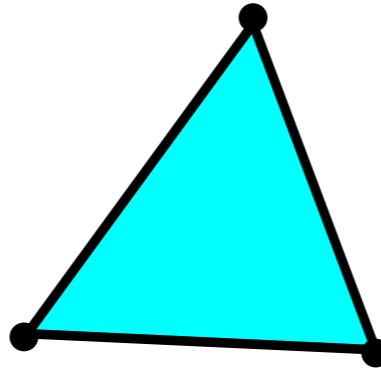
Simplicial complexes



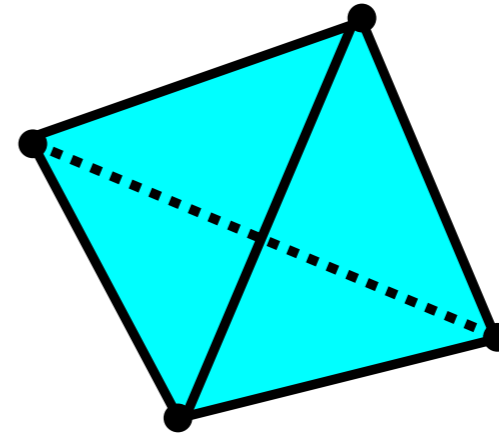
0-simplex:
vertex



1-simplex:
edge



2-simplex:
triangle



3-simplex:
tetrahedron

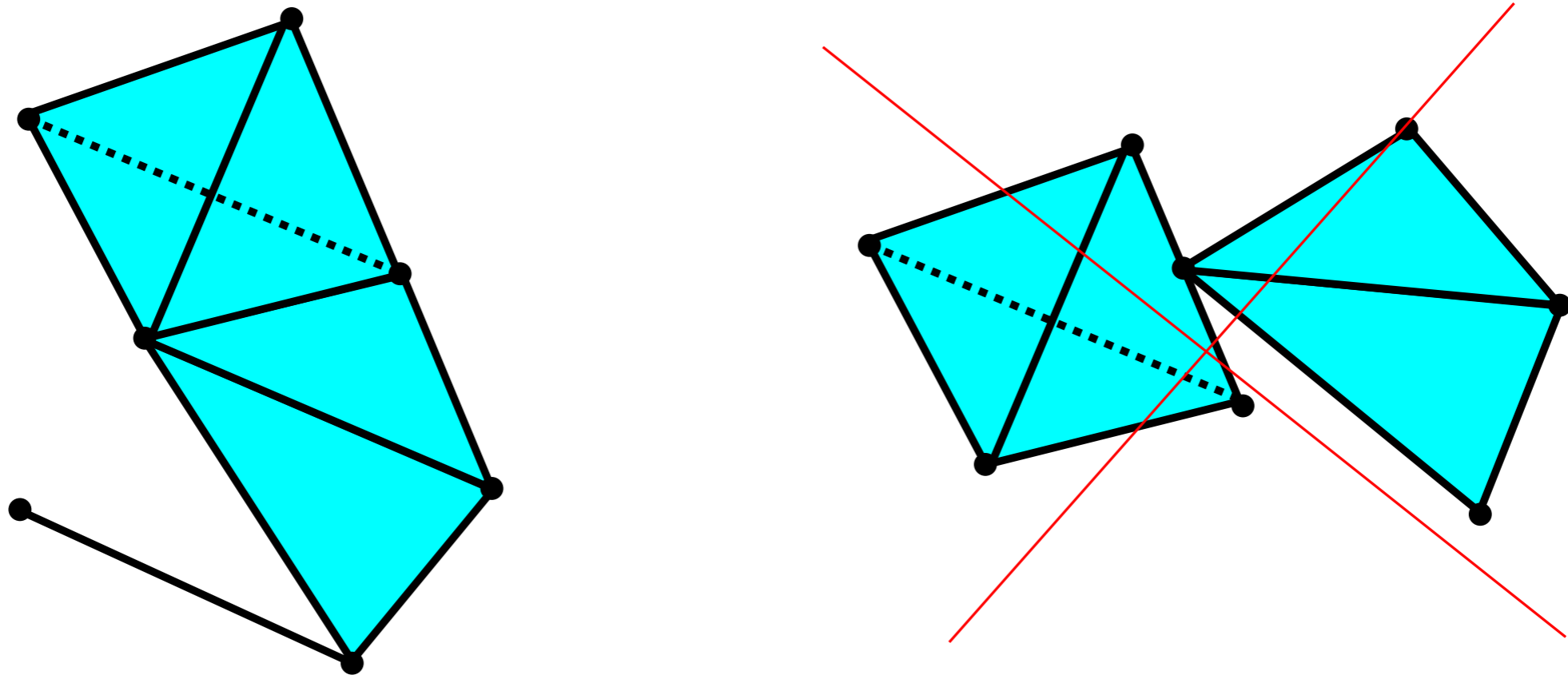
etc...

Given a set $P = \{p_0, \dots, p_k\} \subset \mathbb{R}^d$ of $k + 1$ affinely independent points, the k -dimensional simplex σ , or k -simplex for short, spanned by P is the set of convex combinations

$$\sum_{i=0}^k \lambda_i p_i, \quad \text{with} \quad \sum_{i=0}^k \lambda_i = 1 \quad \text{and} \quad \lambda_i \geq 0.$$

The points p_0, \dots, p_k are called the vertices of σ .

Simplicial complexes



A (finite) **simplicial complex** K in \mathbb{R}^d is a (finite) collection of simplices such that:

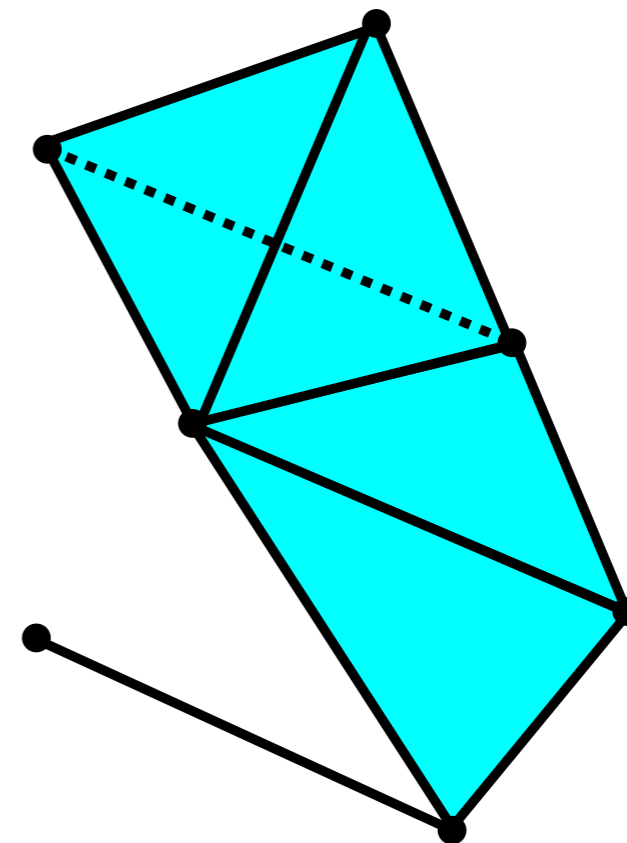
1. any face of a simplex of K is a simplex of K ,
2. the intersection of any two simplices of K is either empty or a common face of both.

The underlying space of K , denoted by $|K| \subset \mathbb{R}^d$ is the union of the simplices of K .

Abstract simplicial complexes

Let $P = \{p_1, \dots, p_n\}$ be a (finite) set. An **abstract simplicial complex** K with vertex set P is a set of subsets of P satisfying the two conditions :

1. The elements of P belong to K .
2. If $\tau \in K$ and $\sigma \subseteq \tau$, then $\sigma \in K$.



The elements of K are the **simplices**.

Let $\{e_1, \dots, e_n\}$ a basis of \mathbb{R}^n . “The” **geometric realization** of K is the (geometric) subcomplex $|K|$ of the simplex spanned by e_1, \dots, e_n such that:

$$[e_{i_0} \cdots e_{i_k}] \in |K| \text{ iff } \{p_{i_0}, \dots, p_{i_k}\} \in K$$

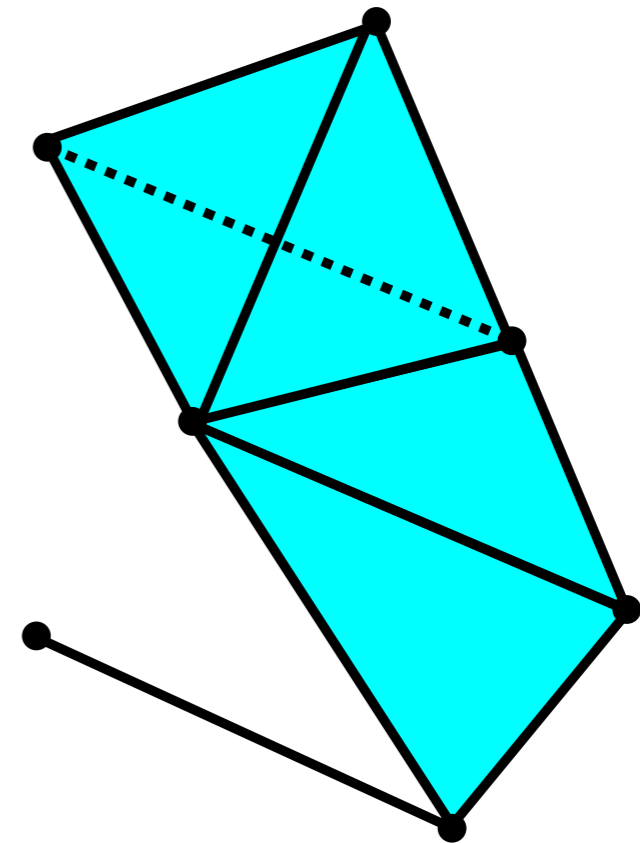
$|K|$ is a topological space (subspace of an Euclidean space)!

Abstract simplicial complexes

Let $P = \{p_1, \dots, p_n\}$ be a (finite) set. An **abstract simplicial complex** K with vertex set P is a set of subsets of P satisfying the two conditions :

1. The elements of P belong to K .
2. If $\tau \in K$ and $\sigma \subseteq \tau$, then $\sigma \in K$.

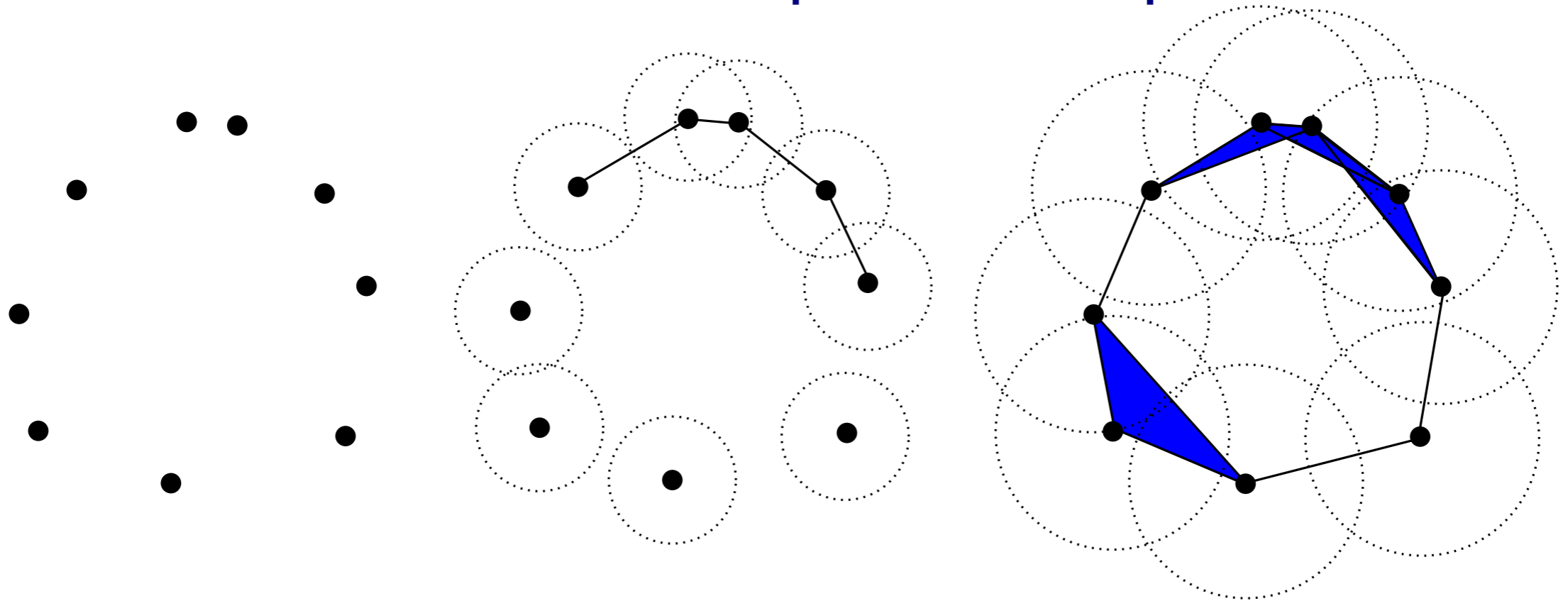
The elements of K are the **simplices**.



IMPORTANT

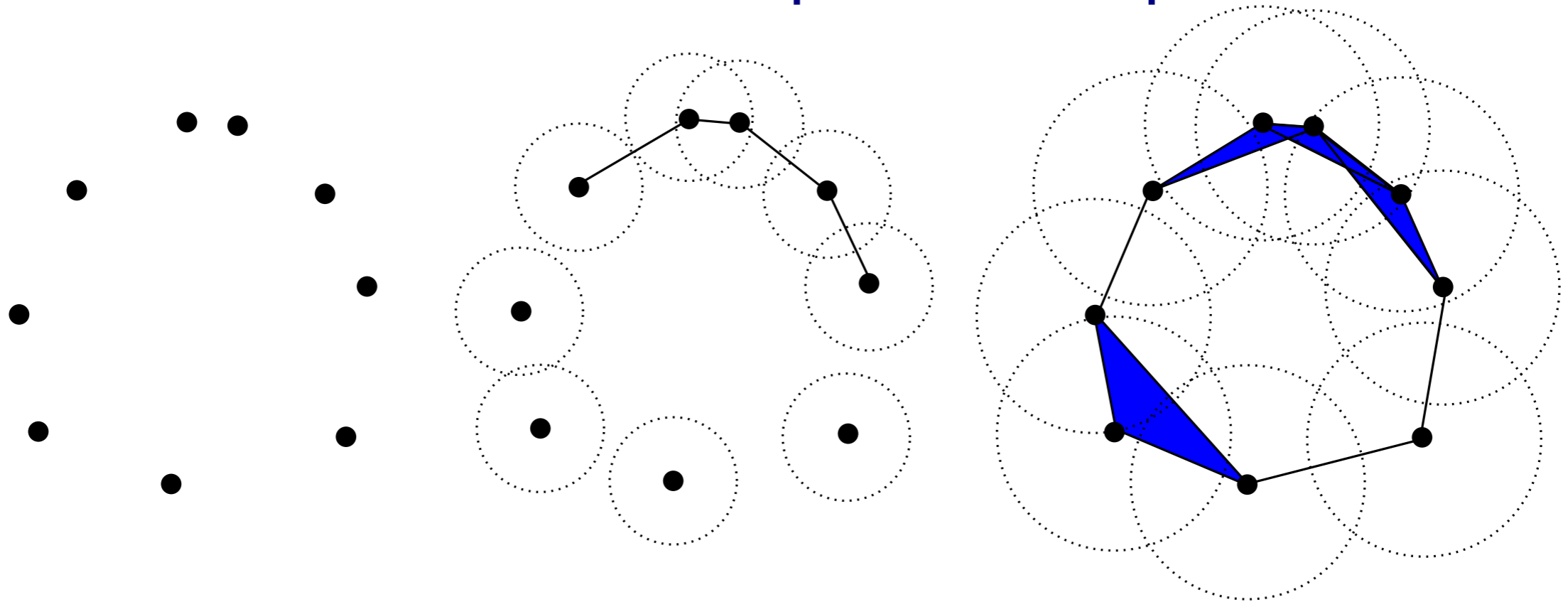
Simplicial complexes can be seen at the same time as geometric/topological spaces (good for top./geom. inference) and as combinatorial objects (abstract simplicial complexes, good for computations).

Filtrations of simplicial complexes



- A **filtered simplicial complex (or a filtration)** \mathbb{S} built on top of a set \mathbb{X} is a family $(\mathbb{S}_a \mid a \in \mathbf{R})$ of subcomplexes of some fixed simplicial complex $\bar{\mathbb{S}}$ with vertex set X s. t. $\mathbb{S}_a \subseteq \mathbb{S}_b$ for any $a \leq b$.
- More generally, **filtration** = nested family of spaces.

Filtrations of simplicial complexes



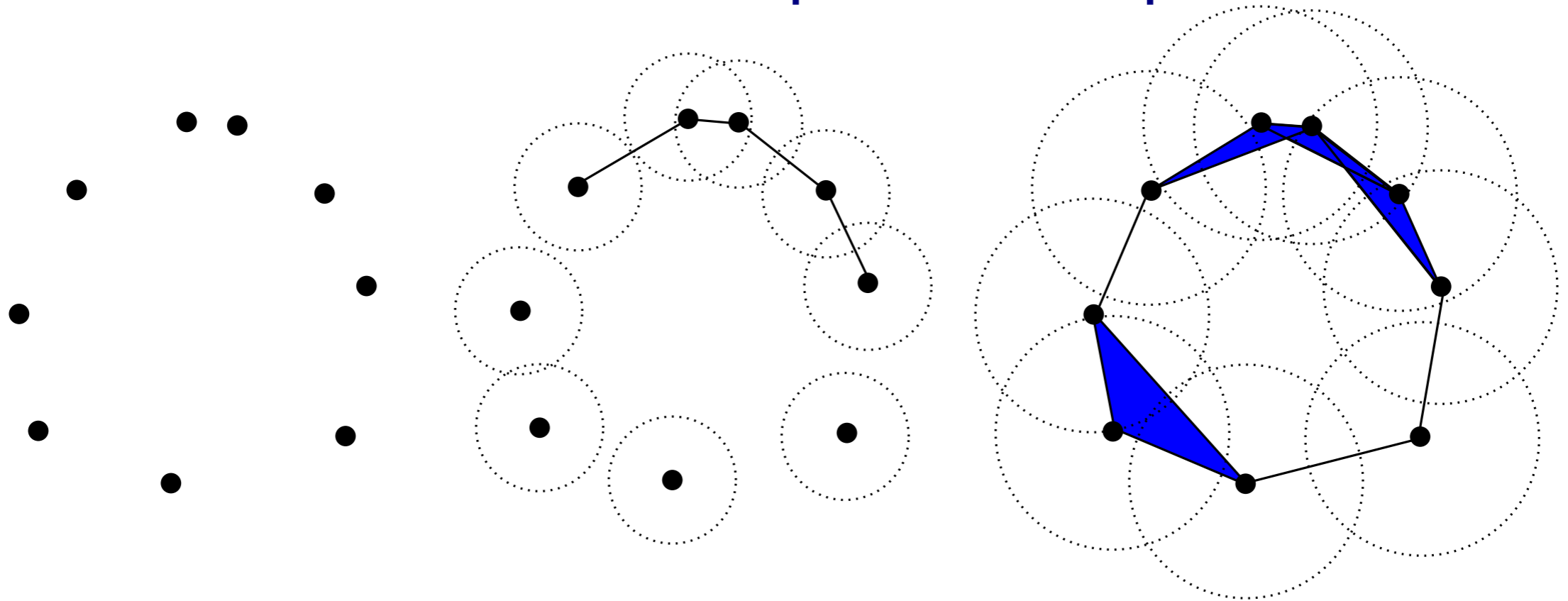
- A **filtered simplicial complex (or a filtration)** \mathbb{S} built on top of a set \mathbb{X} is a family $(\mathbb{S}_a \mid a \in \mathbf{R})$ of subcomplexes of some fixed simplicial complex $\bar{\mathbb{S}}$ with vertex set X s. t. $\mathbb{S}_a \subseteq \mathbb{S}_b$ for any $a \leq b$.
- More generally, **filtration** = nested family of spaces.

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

- The **Vietoris-Rips** filtration is the filtered simplicial complex defined by: for $a \in \mathbf{R}$,

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

Filtrations of simplicial complexes



- A **filtered simplicial complex (or a filtration)** \mathbb{S} built on top of a set \mathbb{X} is a family $(\mathbb{S}_a \mid a \in \mathbf{R})$ of subcomplexes of some fixed simplicial complex $\bar{\mathbb{S}}$ with vertex set X s. t. $\mathbb{S}_a \subseteq \mathbb{S}_b$ for any $a \leq b$.
- More generally, **filtration** = nested family of spaces.

Many other examples and ways to design filtrations depending on the application and targeted objectives : sublevel and upperlevel sets, Čech complex,...

Stability properties

“Stability theorem”: Close spaces/data sets have close persistence diagrams!

[C., de Silva, Oudot - Geom. Dedicata 2013].

If \mathbb{X} and \mathbb{Y} are pre-compact metric spaces, then

$$d_b(\text{dgm}(\text{Rips}(\mathbb{X})), \text{dgm}(\text{Rips}(\mathbb{Y}))) \leq d_{GH}(\mathbb{X}, \mathbb{Y}).$$

Bottleneck distance

Gromov-Hausdorff distance

$$d_{GH}(\mathbb{X}, \mathbb{Y}) := \inf_{\mathbb{Z}, \gamma_1, \gamma_2} d_H(\gamma_1(\mathbb{X}), \gamma_2(\mathbb{Y}))$$

\mathbb{Z} metric space, $\gamma_1 : \mathbb{X} \rightarrow \mathbb{Z}$ and $\gamma_2 : \mathbb{Y} \rightarrow \mathbb{Z}$
isometric embeddings.

Rem: This result also holds for other families of filtrations (particular case of a more general thm).

Stability properties

“Stability theorem”: Close spaces/data sets have close persistence diagrams!

[C., de Silva, Oudot - Geom. Dedicata 2013].

If \mathbb{X} and \mathbb{Y} are pre-compact metric spaces, then

$$d_b(\text{dgm}(\text{Rips}(\mathbb{X})), \text{dgm}(\text{Rips}(\mathbb{Y}))) \leq d_{GH}(\mathbb{X}, \mathbb{Y}).$$

Bottleneck distance

Gromov-Hausdorff distance

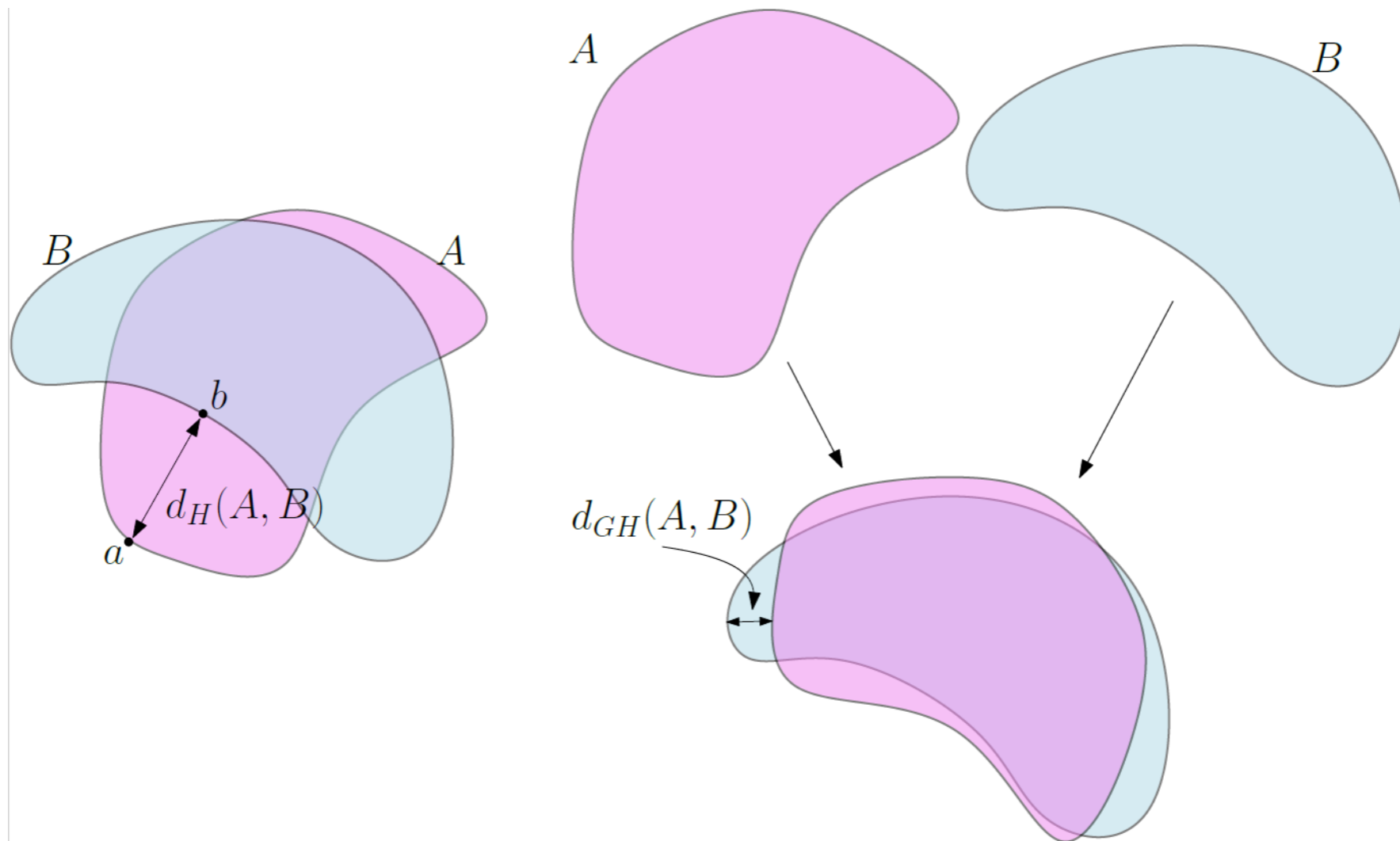
$$d_{GH}(\mathbb{X}, \mathbb{Y}) := \inf_{\mathbb{Z}, \gamma_1, \gamma_2} d_H(\gamma_1(\mathbb{X}), \gamma_2(\mathbb{Y}))$$

\mathbb{Z} metric space, $\gamma_1 : \mathbb{X} \rightarrow \mathbb{Z}$ and $\gamma_2 : \mathbb{Y} \rightarrow \mathbb{Z}$
isometric embeddings.

Rem: This result also holds for other families of filtrations (particular case of a more general thm).

From a statistical perspective, when \mathbb{X} is a random point cloud, such result links the study of statistical properties of persistence diagrams to support estimation problems.

Hausdorff distance



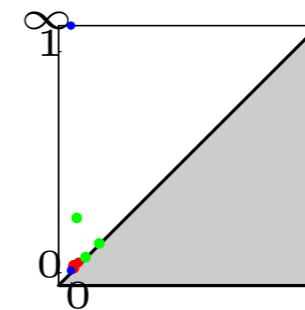
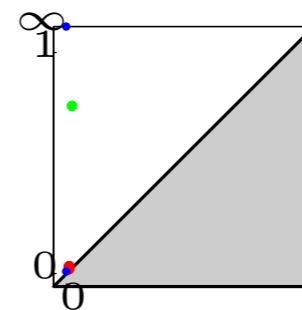
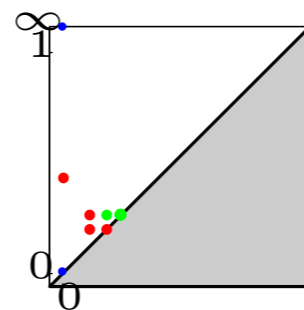
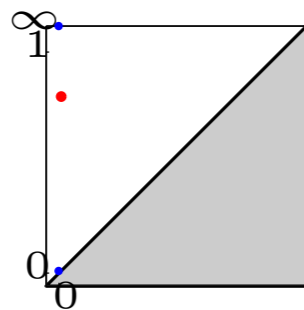
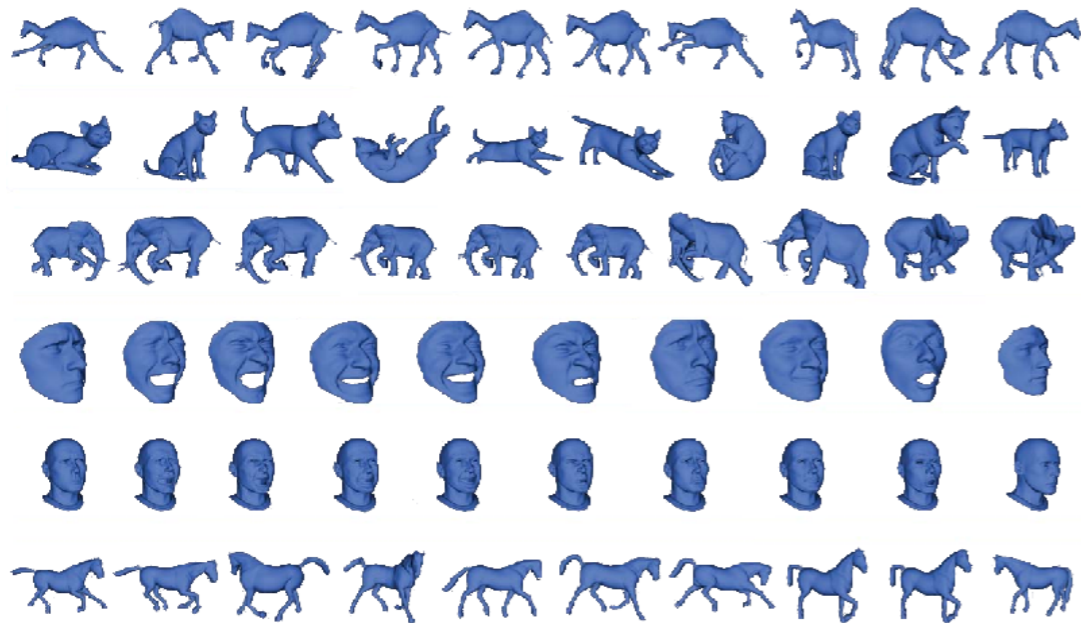
Let $A, B \subset M$ be two compact subsets of a metric space (M, d)

$$d_H(A, B) = \max\left\{\sup_{b \in B} d(b, A), \sup_{a \in A} d(a, B)\right\}$$

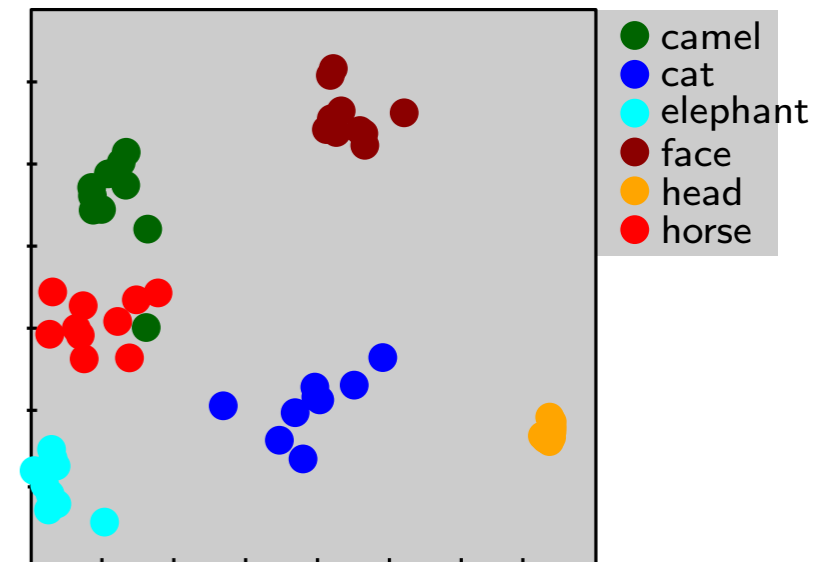
where $d(b, A) = \sup_{a \in A} d(b, a)$.

Application: non rigid shape classification

[C., Cohen-Steiner, Guibas, Mémoli, Oudot - SGP '09]



MDS using bottleneck distance.



- Non rigid shapes in a same class are almost isometric, but computing Gromov-Hausdorff distance between shapes is extremely expensive.
- Compare diagrams of sampled shapes instead of shapes themselves.

Where do stability results come from?

Definition: A **persistence module** \mathbb{V} is an indexed family of vector spaces $(V_a \mid a \in \mathbf{R})$ and a doubly-indexed family of linear maps $(v_a^b : V_a \rightarrow V_b \mid a \leq b)$ which satisfy the composition law $v_b^c \circ v_a^b = v_a^c$ whenever $a \leq b \leq c$, and where v_a^a is the identity map on V_a .

Examples:

- Let \mathbb{S} be a filtered simplicial complex. If $V_a = H(\mathbb{S}_a)$ and $v_a^b : H(\mathbb{S}_a) \rightarrow H(\mathbb{S}_b)$ is the linear map induced by the inclusion $\mathbb{S}_a \hookrightarrow \mathbb{S}_b$ then $(H(\mathbb{S}_a) \mid a \in \mathbf{R})$ is a persistence module.
- Given a metric space $(\mathbb{X}, d_{\mathbb{X}})$, $H(\text{Rips}(\mathbb{X}))$ is a persistence module.
- If $f : X \rightarrow \mathbf{R}$ is a function, then the filtration defined by the sublevel sets of f , $\mathbb{F}_a = f^{-1}((-\infty, a])$, induces a persistence module at homology level.

Where do stability results come from?

Definition: A **persistence module** \mathbb{V} is an indexed family of vector spaces $(V_a \mid a \in \mathbf{R})$ and a doubly-indexed family of linear maps $(v_a^b : V_a \rightarrow V_b \mid a \leq b)$ which satisfy the composition law $v_b^c \circ v_a^b = v_a^c$ whenever $a \leq b \leq c$, and where v_a^a is the identity map on V_a .

Definition: A persistence module \mathbb{V} is **q-tame** if for any $a < b$, v_a^b has a finite rank.

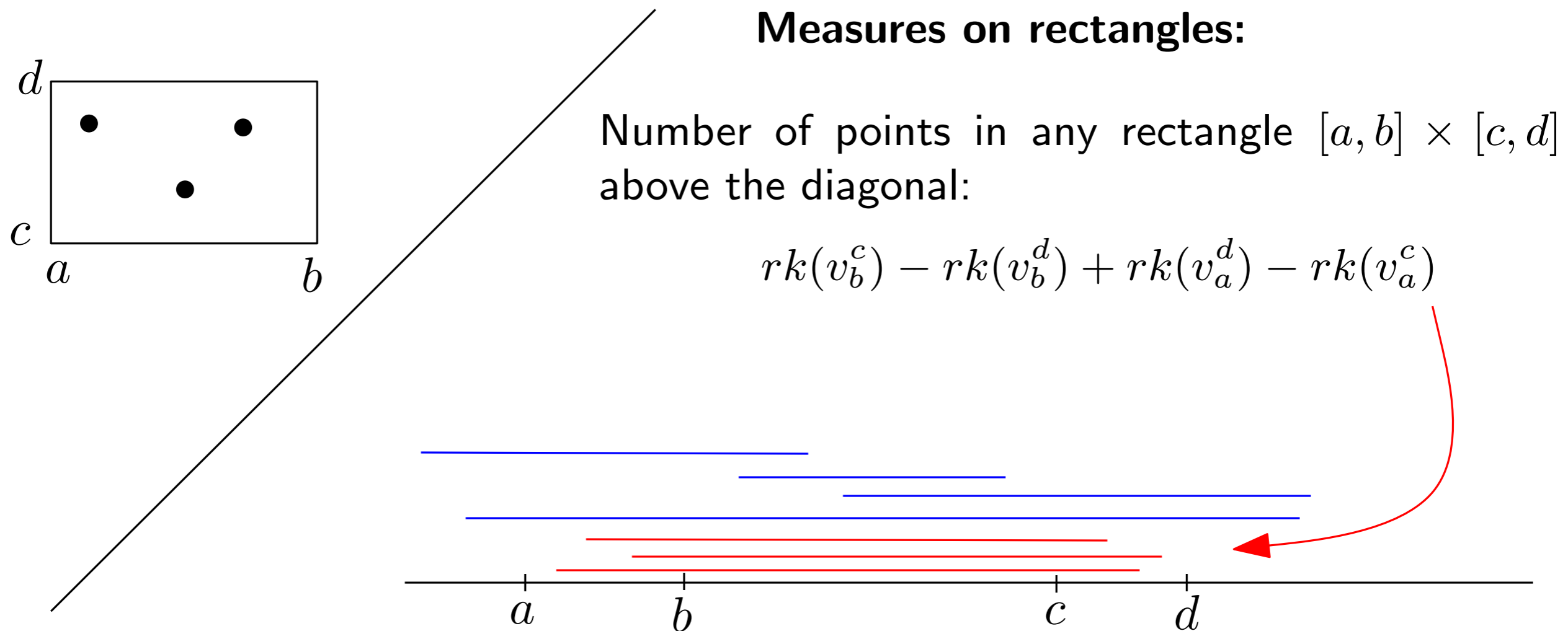
Theorem: [C., Cohen-Steiner, Glisse, Guibas, Oudot - SoCG'09], [C., de Silva, Glisse, Oudot 12]

q-tame persistence modules have well-defined persistence diagrams.

Where do stability results come from?

Definition: A **persistence module** \mathbb{V} is an indexed family of vector spaces $(V_a \mid a \in \mathbf{R})$ and a doubly-indexed family of linear maps $(v_a^b : V_a \rightarrow V_b \mid a \leq b)$ which satisfy the composition law $v_b^c \circ v_a^b = v_a^c$ whenever $a \leq b \leq c$, and where v_a^a is the identity map on V_a .

An idea about the definition of persistence diagrams:



Where do stability results come from?

Definition: A **persistence module** \mathbb{V} is an indexed family of vector spaces $(V_a \mid a \in \mathbf{R})$ and a doubly-indexed family of linear maps $(v_a^b : V_a \rightarrow V_b \mid a \leq b)$ which satisfy the composition law $v_b^c \circ v_a^b = v_a^c$ whenever $a \leq b \leq c$, and where v_a^a is the identity map on V_a .

Definition: A persistence module \mathbb{V} is **q-tame** if for any $a < b$, v_a^b has a finite rank.

Theorem: [C., Cohen-Steiner, Glisse, Guibas, Oudot - SoCG'09], [C., de Silva, Glisse, Oudot 12]

q-tame persistence modules have well-defined persistence diagrams.

Exercise: Let \mathbb{X} be a precompact metric space. Then $H(\text{Rips}(\mathbb{X}))$ is q-tame.

Recall that a metric space (\mathbb{X}, ρ) is **precompact** if for any $\epsilon > 0$ there exists a finite subset $F_\epsilon \subset \mathbb{X}$ such that $d_H(\mathbb{X}, F_\epsilon) < \epsilon$ (i.e. $\forall x \in X, \exists p \in F_\epsilon$ s.t. $\rho(x, p) < \epsilon$).

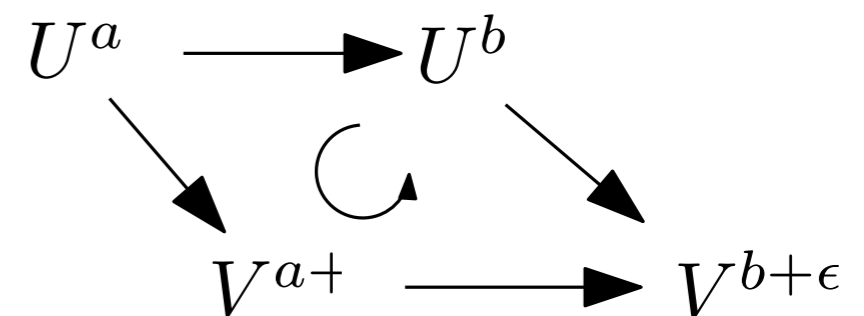
Where do stability results come from?

Definition: A **persistence module** \mathbb{V} is an indexed family of vector spaces $(V_a \mid a \in \mathbf{R})$ and a doubly-indexed family of linear maps $(v_a^b : V_a \rightarrow V_b \mid a \leq b)$ which satisfy the composition law $v_b^c \circ v_a^b = v_a^c$ whenever $a \leq b \leq c$, and where v_a^a is the identity map on V_a .

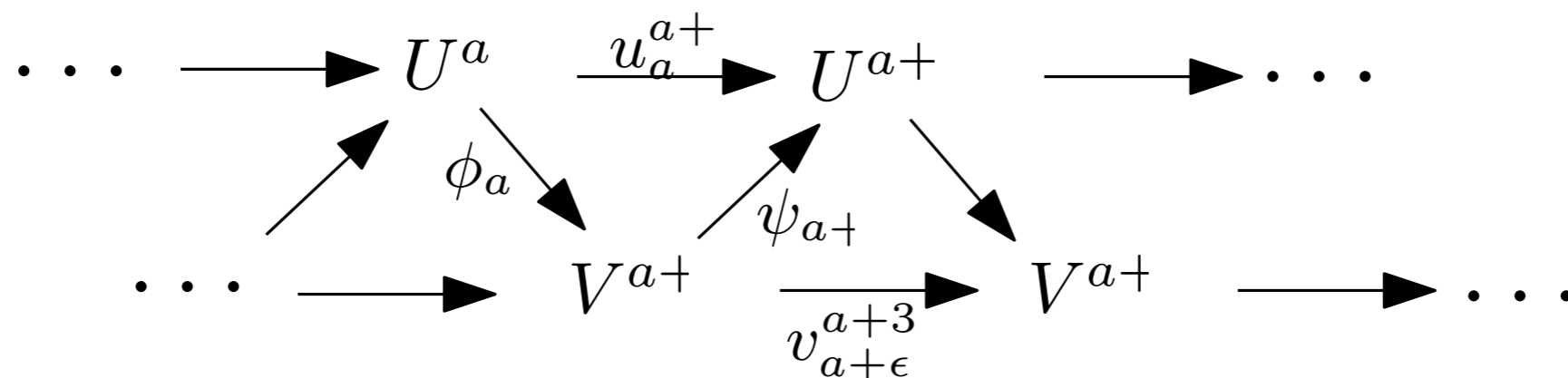
A **homomorphism of degree ϵ** between two persistence modules \mathbb{U} and \mathbb{V} is a collection Φ of linear maps

$$(\phi_a : U_a \rightarrow V_{a+\epsilon} \mid a \in \mathbf{R})$$

such that $v_{a+\epsilon}^{b+\epsilon} \circ \phi_a = \phi_b \circ u_a^b$ for all $a \leq b$.



An **ϵ -interleaving** between \mathbb{U} and \mathbb{V} is specified by two homomorphisms of degree ϵ $\Phi : \mathbb{U} \rightarrow \mathbb{V}$ and $\Psi : \mathbb{V} \rightarrow \mathbb{U}$ s.t. $\Phi \circ \Psi$ and $\Psi \circ \Phi$ are the “shifts” of degree 2ϵ between \mathbb{U} and \mathbb{V} .



Where do stability results come from?

Definition: A **persistence module** \mathbb{V} is an indexed family of vector spaces $(V_a \mid a \in \mathbf{R})$ and a doubly-indexed family of linear maps $(v_a^b : V_a \rightarrow V_b \mid a \leq b)$ which satisfy the composition law $v_b^c \circ v_a^b = v_a^c$ whenever $a \leq b \leq c$, and where v_a^a is the identity map on V_a .

Stability Thm: [C., Cohen-Steiner, Glisse, Guibas, Oudot - SoCG '09], [C., de Silva, Glisse, Oudot 12]

If \mathbb{U} and \mathbb{V} are q -tame and ϵ -interleaved for some $\epsilon \geq 0$ then

$$d_B(\mathrm{dgm}(\mathbb{U}), \mathrm{dgm}(\mathbb{V})) \leq \epsilon$$

Where do stability results come from?

Definition: A **persistence module** \mathbb{V} is an indexed family of vector spaces $(V_a \mid a \in \mathbf{R})$ and a doubly-indexed family of linear maps $(v_a^b : V_a \rightarrow V_b \mid a \leq b)$ which satisfy the composition law $v_b^c \circ v_a^b = v_a^c$ whenever $a \leq b \leq c$, and where v_a^a is the identity map on V_a .

Stability Thm: [C., Cohen-Steiner, Glisse, Guibas, Oudot - SoCG '09], [C., de Silva, Glisse, Oudot 12]

If \mathbb{U} and \mathbb{V} are q -tame and ϵ -interleaved for some $\epsilon \geq 0$ then

$$d_B(\mathrm{dgm}(\mathbb{U}), \mathrm{dgm}(\mathbb{V})) \leq \epsilon$$

Strategy: build filtrations that induce **q-tame** homology persistence modules and that turn out to be **ϵ -interleaved** when the considered spaces/functions are $O(\epsilon)$ -close.

Why persistence

- Even when X is compact, $H_p(\text{Rips}(X, a))$, $p \geq 1$, might be infinite dimensional for some value of a :

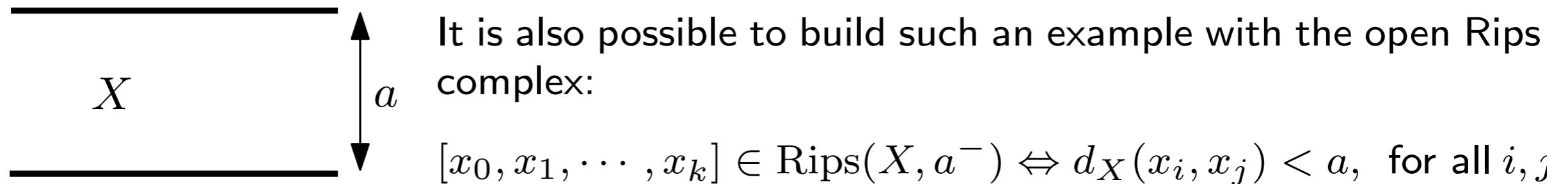


It is also possible to build such an example with the open Rips complex:

$$[x_0, x_1, \dots, x_k] \in \text{Rips}(X, a^-) \Leftrightarrow d_X(x_i, x_j) < a, \text{ for all } i, j$$

Why persistence

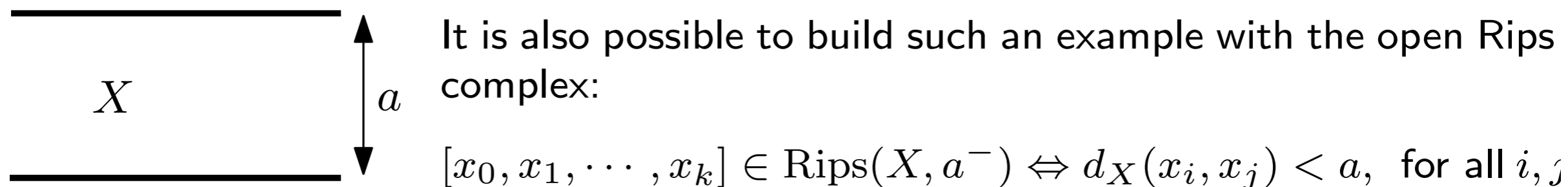
- Even when X is compact, $H_p(\text{Rips}(X, a))$, $p \geq 1$, might be infinite dimensional for some value of a :



- For any $\alpha, \beta \in \mathbf{R}$ such that $0 < \alpha \leq \beta$ and any integer k there exists a compact metric space X such that for any $a \in [\alpha, \beta]$, $H_k(\text{Rips}(X, a))$ has a non countable infinite dimension (can be embedded in \mathbf{R}^4 [Droz 2013]).

Why persistence

- Even when X is compact, $H_p(\text{Rips}(X, a))$, $p \geq 1$, might be infinite dimensional for some value of a :



- For any $\alpha, \beta \in \mathbf{R}$ such that $0 < \alpha \leq \beta$ and any integer k there exists a compact metric space X such that for any $a \in [\alpha, \beta]$, $H_k(\text{Rips}(X, a))$ has a non countable infinite dimension (can be embedded in \mathbf{R}^4 [Droz 2013]).
- If X is compact, then $\dim H_1(\check{\text{Cech}}(X, a)) < +\infty$ for all a ([Smale-Smale, C.-de Silva]).
- If X is geodesic, then $\dim H_1(\text{Rips}(X, a)) < +\infty$ for all $a > 0$ and $\text{Dgm}(H_1(\mathbb{R}\text{ips}(X)))$ is contained in the vertical line $x = 0$.
- If X is a geodesic δ -hyperbolic space then $\text{Dgm}(H_2(\mathbb{R}\text{ips}(X)))$ is contained in a vertical band of width $O(\delta)$.

Persistent homology with the GUDHI library



गुढी **GUDHI** Geometry Understanding
in Higher Dimensions

<http://gudhi.gforge.inria.fr/>

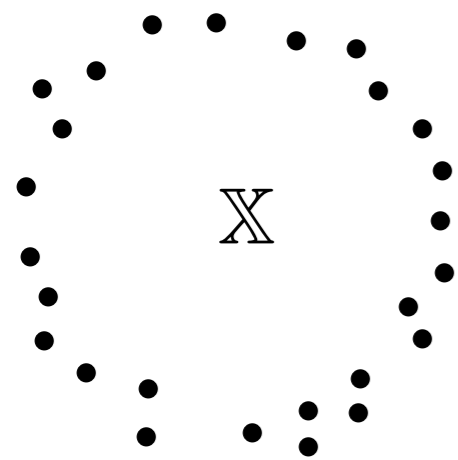
GUDHI :

- a C++/Python open source software library for TDA,
- a developers team, an editorial board, open to external contributions,
- provides state-of-the-art TDA data structures and algorithms : design of filtrations, computation of pre-defined filtrations, persistence diagrams,...
- part of GUDHI is interfaced to R through the TDA package.

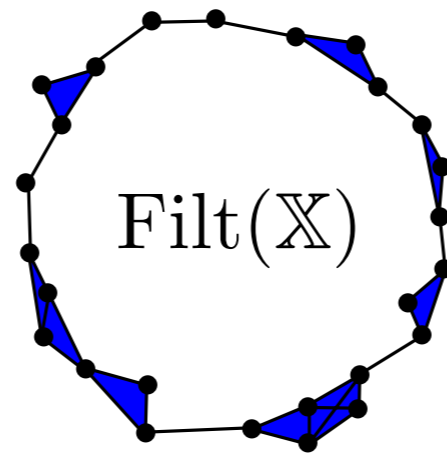
Statistical properties and features extraction from
persistence diagrams

Statistical setting and “linear representations”

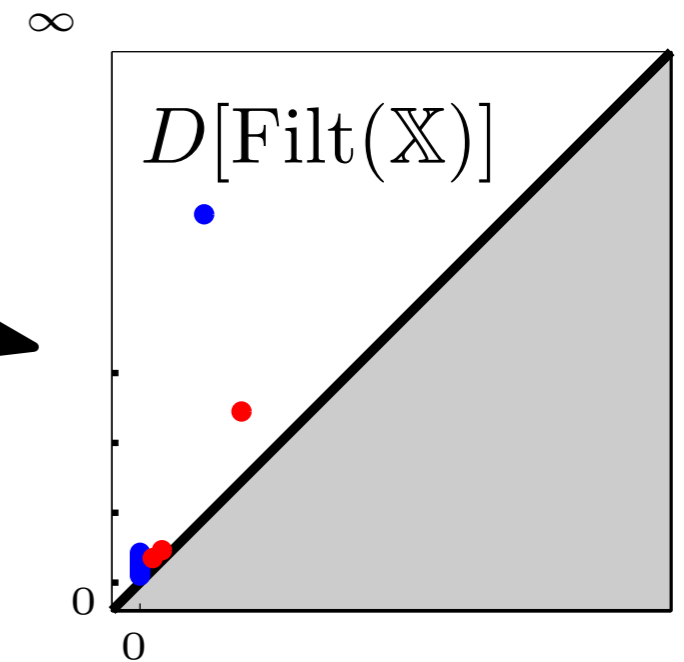
\mathbb{X} is now a random point cloud (in some metric space)



Filt is a deterministic filtration (e.g. Rips)



$D[\text{Filt}(\mathbb{X})]$ becomes random

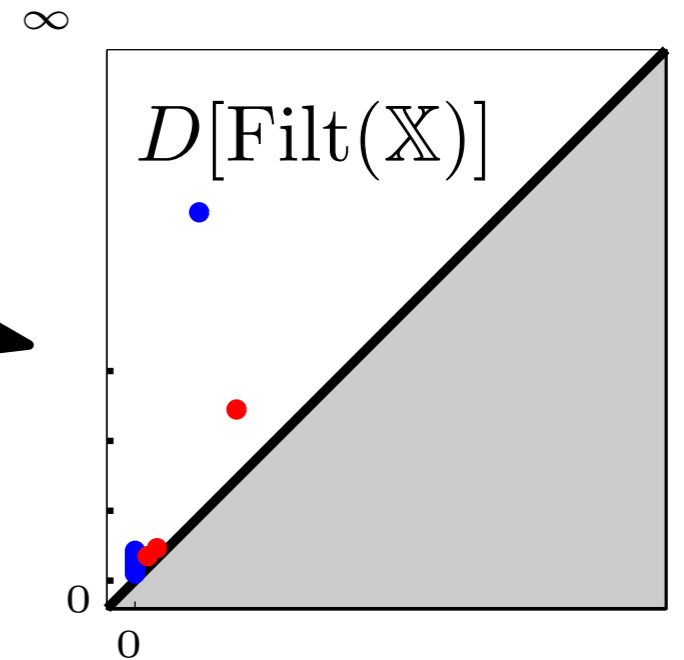
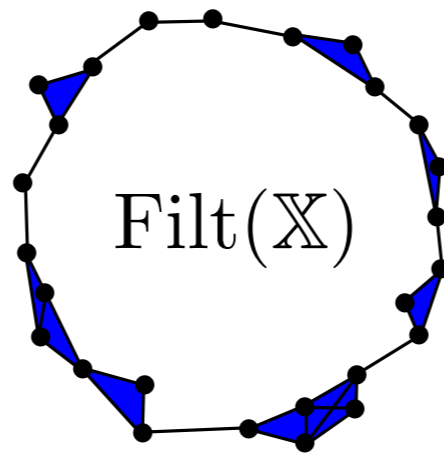
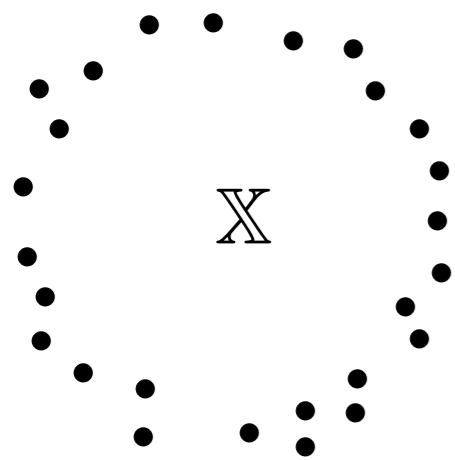


Statistical setting and “linear representations”

\mathbb{X} is now a random point cloud (in some metric space)

Filt is a deterministic filtration (e.g. Rips)

$D[\text{Filt}(\mathbb{X})]$ becomes random



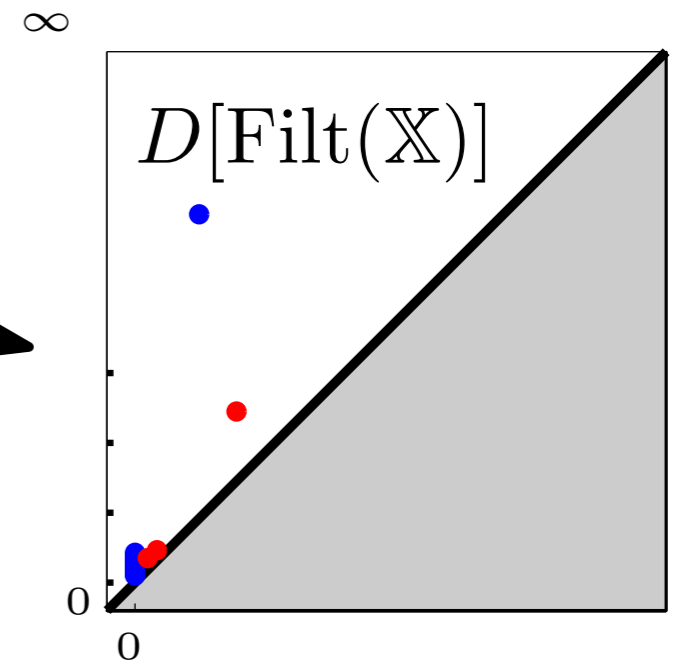
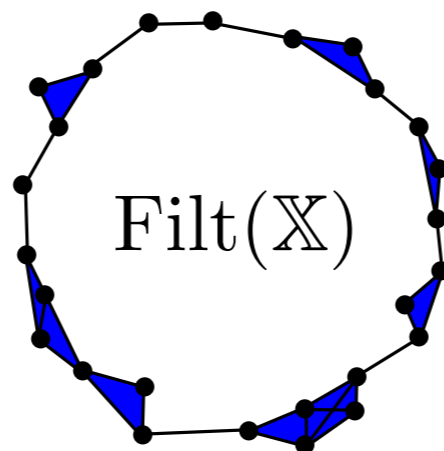
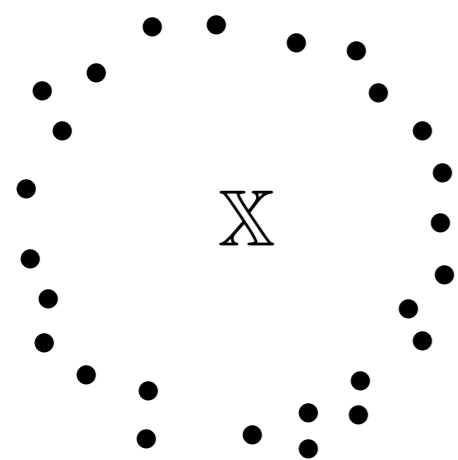
What can be said about the distribution of diagrams $D[\text{Filt}(\mathbb{X})]$?

Statistical setting and “linear representations”

\mathbb{X} is now a random point cloud (in some metric space)

Filt is a deterministic filtration (e.g. Rips)

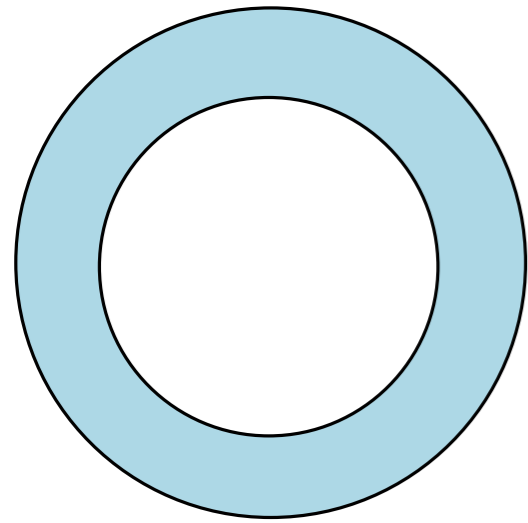
$D[\text{Filt}(\mathbb{X})]$ becomes random



What can be said about the distribution of diagrams $D[\text{Filt}(\mathbb{X})]$?

- Stability properties \Rightarrow asymptotic properties, confidence bands, Wasserstein stability,...
- Other representation of persistence that are well-suited for ML (landscapes, Betti curves, pers. images, kernels,...)

Statistical setting



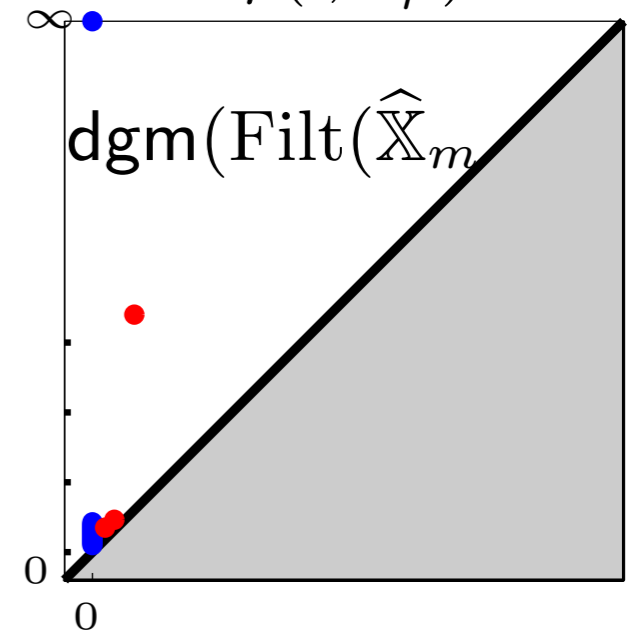
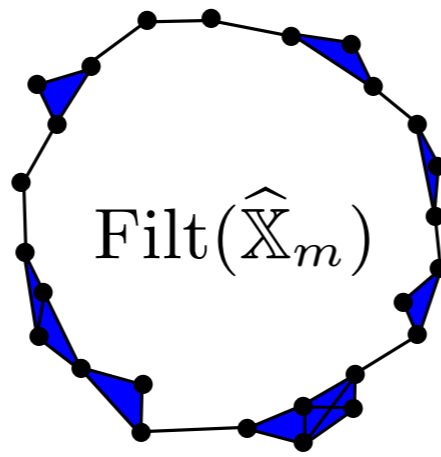
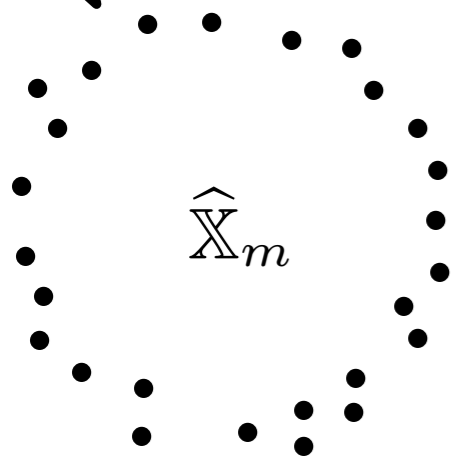
(\mathbb{M}, ρ) metric space

μ a probability measure with **compact** support \mathbb{X}_μ .

Sample m points
according to μ .

Examples:

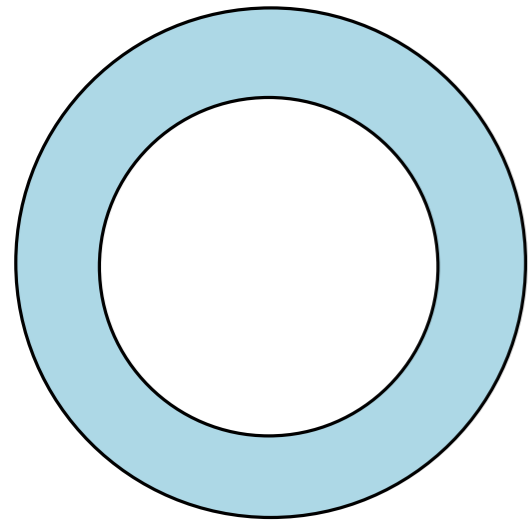
- $\text{Filt}(\hat{\mathbb{X}}_m) = \text{Rips}_\alpha(\hat{\mathbb{X}}_m)$
- $\text{Filt}(\hat{\mathbb{X}}_m) = \check{\text{Cech}}_\alpha(\hat{\mathbb{X}}_m)$
- $\text{Filt}(\hat{\mathbb{X}}_m) = \text{sublevelset filtration of } \rho(., \mathbb{X}_\mu).$



Questions:

- Statistical properties of $\text{dgm}(\text{Filt}(\hat{\mathbb{X}}_m))$? $\text{dgm}(\text{Filt}(\hat{\mathbb{X}}_m)) \rightarrow ?$ as $m \rightarrow +\infty$?

Statistical setting



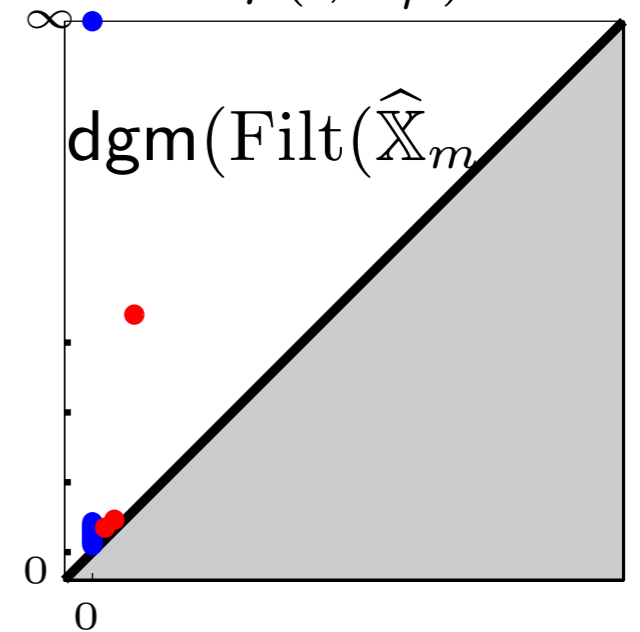
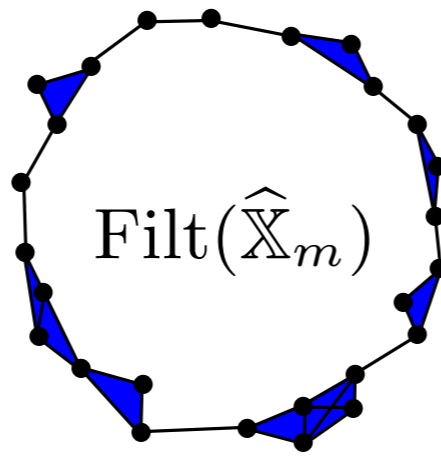
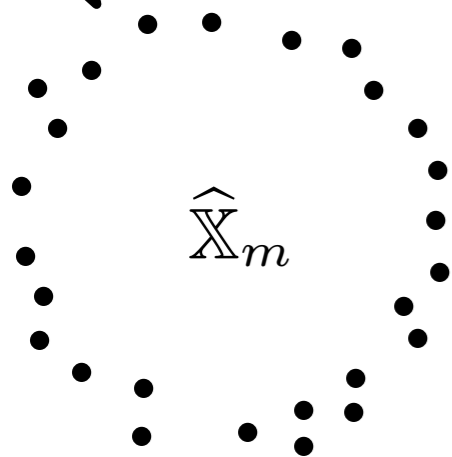
(\mathbb{M}, ρ) metric space

μ a probability measure with **compact** support \mathbb{X}_μ .

Sample m points
according to μ .

Examples:

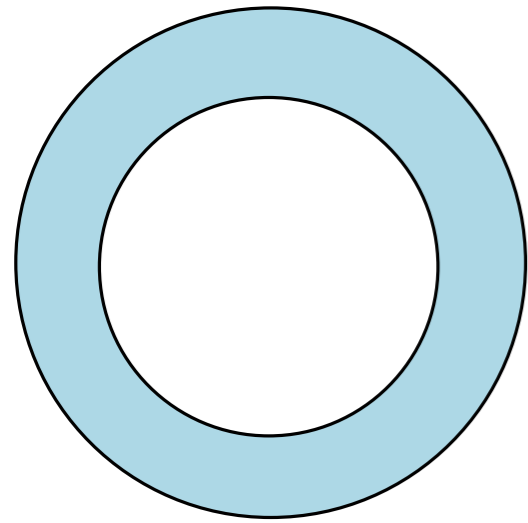
- $\text{Filt}(\hat{\mathbb{X}}_m) = \text{Rips}_\alpha(\hat{\mathbb{X}}_m)$
- $\text{Filt}(\hat{\mathbb{X}}_m) = \check{\text{Cech}}_\alpha(\hat{\mathbb{X}}_m)$
- $\text{Filt}(\hat{\mathbb{X}}_m) = \text{sublevelset filtration of } \rho(., \mathbb{X}_\mu).$



Questions:

- Statistical properties of $\text{dgm}(\text{Filt}(\hat{\mathbb{X}}_m))$? $\text{dgm}(\text{Filt}(\hat{\mathbb{X}}_m)) \rightarrow ?$ as $m \rightarrow +\infty$?
- Can we do more statistics with persistence diagrams? What can be said about distributions of diagrams?

Statistical setting



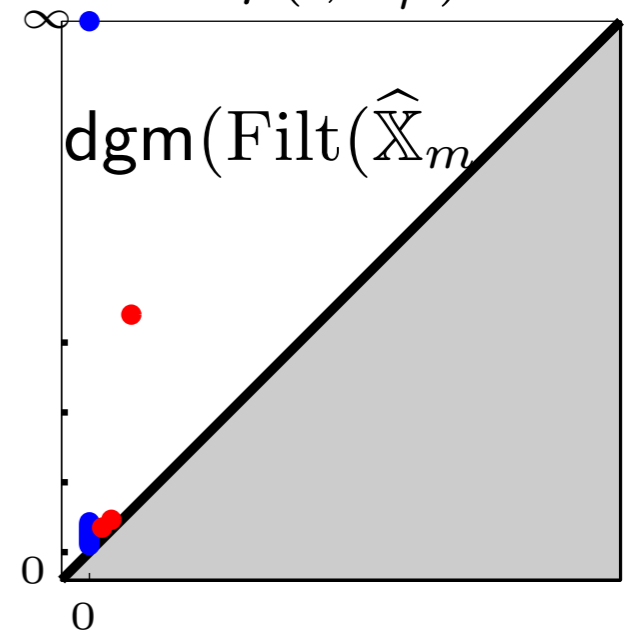
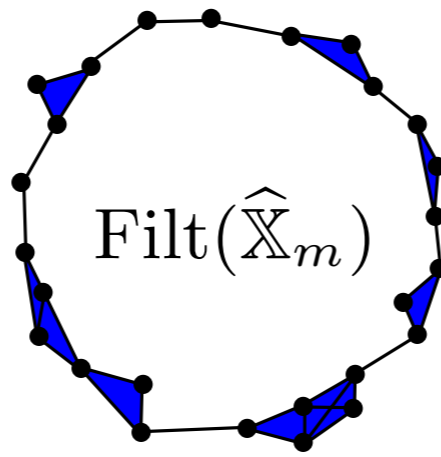
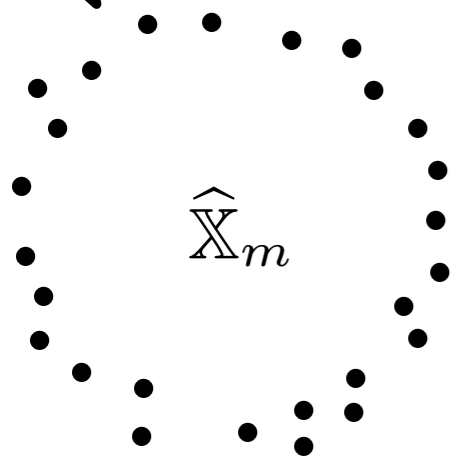
(\mathbb{M}, ρ) metric space

μ a probability measure with **compact** support \mathbb{X}_μ .

Sample m points
according to μ .

Examples:

- $\text{Filt}(\hat{\mathbb{X}}_m) = \text{Rips}_\alpha(\hat{\mathbb{X}}_m)$
- $\text{Filt}(\hat{\mathbb{X}}_m) = \check{\text{Cech}}_\alpha(\hat{\mathbb{X}}_m)$
- $\text{Filt}(\hat{\mathbb{X}}_m) = \text{sublevelset filtration of } \rho(\cdot, \mathbb{X}_\mu).$



Stability thm: $d_b(\text{dgm}(\text{Filt}(\mathbb{X}_\mu)), \text{dgm}(\text{Filt}(\hat{\mathbb{X}}_m))) \leq 2d_{GH}(\mathbb{X}_\mu, \hat{\mathbb{X}}_m)$

So, for any $\varepsilon > 0$,

$$\mathbb{P} \left(d_b \left(\text{dgm}(\text{Filt}(\mathbb{X}_\mu)), \text{dgm}(\text{Filt}(\hat{\mathbb{X}}_m)) \right) > \varepsilon \right) \leq \mathbb{P} \left(d_{GH}(\mathbb{X}_\mu, \hat{\mathbb{X}}_m) > \frac{\varepsilon}{2} \right)$$

Deviation inequality and rate of convergence

[C., Glisse, Labruère, Michel ICML'14 - JMLR'15]

For $a, b > 0$, μ satisfies the (a, b) -standard assumption if for any $x \in \mathbb{X}_\mu$ and any $r > 0$, we have $\mu(B(x, r)) \geq \min(ar^b, 1)$.

Deviation inequality and rate of convergence

[C., Glisse, Labruère, Michel ICML'14 - JMLR'15]

For $a, b > 0$, μ satisfies the (a, b) -standard assumption if for any $x \in \mathbb{X}_\mu$ and any $r > 0$, we have $\mu(B(x, r)) \geq \min(ar^b, 1)$.

Theorem: If μ satisfies the (a, b) -standard assumption, then for any $\varepsilon > 0$:

$$\mathbb{P} \left(d_b \left(\text{dgm}(\text{Filt}(\mathbb{X}_\mu)), \text{dgm}(\text{Filt}(\widehat{\mathbb{X}}_m)) \right) > \varepsilon \right) \leq \min\left(\frac{8^b}{a\varepsilon^b} \exp(-ma\varepsilon^b), 1\right).$$

Deviation inequality and rate of convergence

[C., Glisse, Labruère, Michel ICML'14 - JMLR'15]

For $a, b > 0$, μ satisfies the (a, b) -standard assumption if for any $x \in \mathbb{X}_\mu$ and any $r > 0$, we have $\mu(B(x, r)) \geq \min(ar^b, 1)$.

Theorem: If μ satisfies the (a, b) -standard assumption, then for any $\varepsilon > 0$:

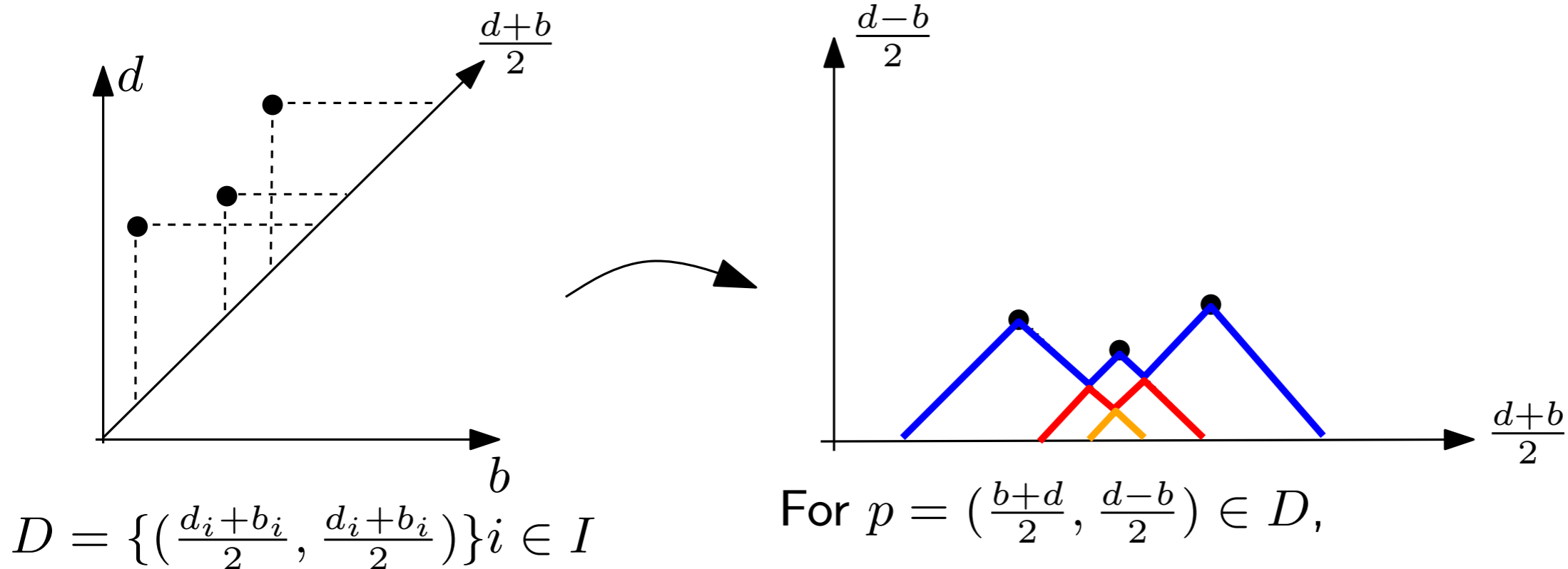
$$\mathbb{P} \left(d_b \left(\text{dgm}(\text{Filt}(\mathbb{X}_\mu)), \text{dgm}(\text{Filt}(\hat{\mathbb{X}}_m)) \right) > \varepsilon \right) \leq \min\left(\frac{8^b}{a\varepsilon^b} \exp(-ma\varepsilon^b), 1\right).$$

Corollary: Let $\mathcal{P}(a, b, \mathbb{M})$ be the set of (a, b) -standard proba measures on \mathbb{M} . Then:

$$\sup_{\mu \in \mathcal{P}(a, b, \mathbb{M})} \mathbb{E} \left[d_b(\text{dgm}(\text{Filt}(\mathbb{X}_\mu)), \text{dgm}(\text{Filt}(\hat{\mathbb{X}}_m))) \right] \leq C \left(\frac{\ln m}{m} \right)^{1/b}$$

where the constant C only depends on a and b (**not on \mathbb{M} !**). Moreover, **the upper bound is tight (in a minimax sense)!**

Persistence landscapes



$$\Lambda_p(t) = \begin{cases} t - b & t \in [b, \frac{b+d}{2}] \\ d - t & t \in (\frac{b+d}{2}, d] \\ 0 & \text{otherwise.} \end{cases}$$

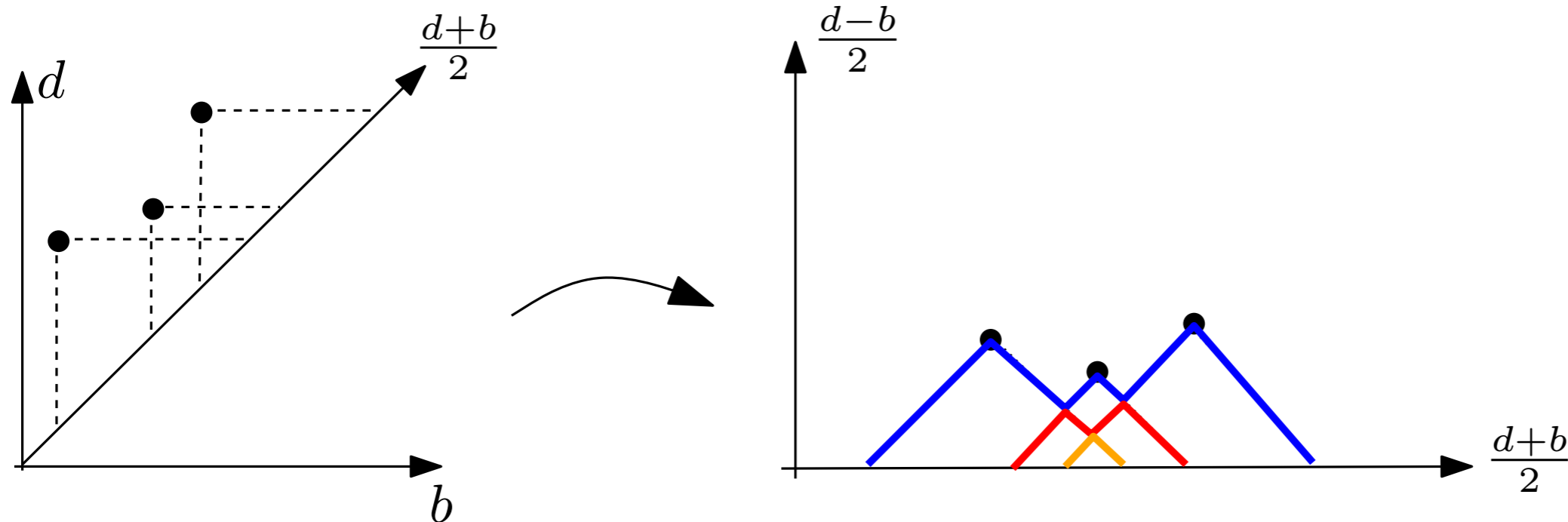
Persistence landscape [Bubenik 2012]:

$$\lambda_D(k, t) = \text{kmax}_{p \in \text{dgm}} \Lambda_p(t), \quad t \in \mathbb{R}, k \in \mathbb{N},$$

where kmax is the k th largest value in the set.

Many other ways to “linearize” persistence diagrams: intensity functions, image persistence, Betti curves, kernels,...

Persistence landscapes



Persistence landscape [Bubenik 2012]:

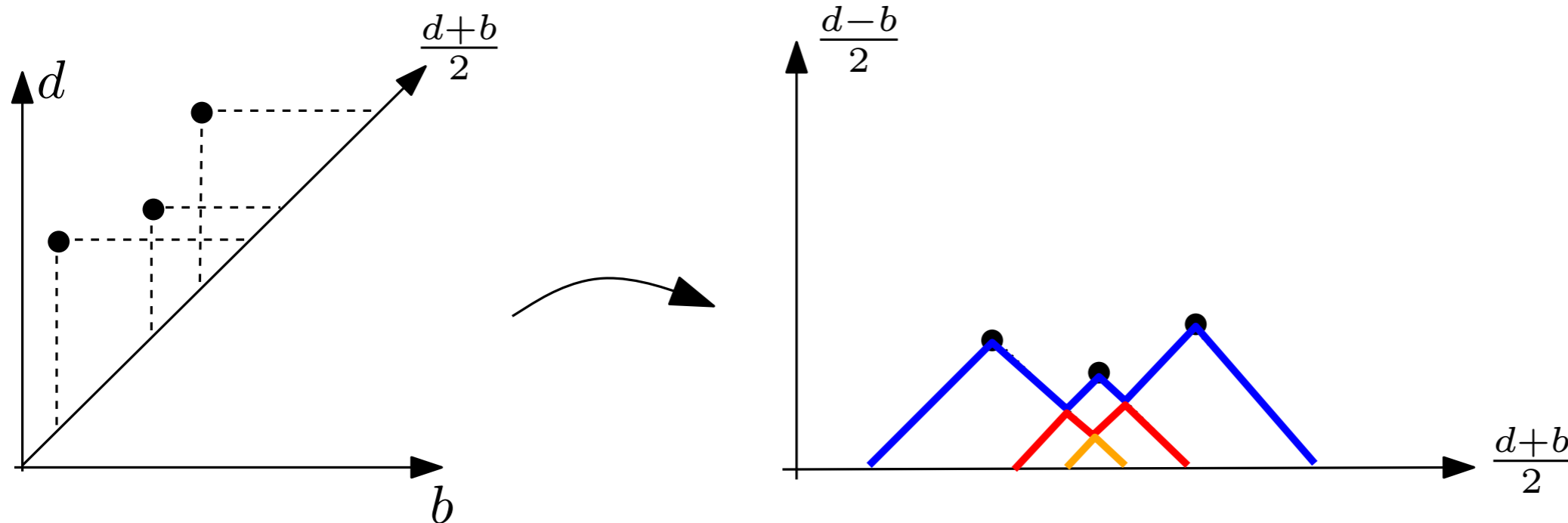
$$\lambda_D(k, t) = k \max_{p \in \text{dgm}} \Lambda_p(t), \quad t \in \mathbb{R}, k \in \mathbb{N},$$

Properties

- For any $t \in \mathbb{R}$ and any $k \in \mathbb{N}$, $0 \leq \lambda_D(k, t) \leq \lambda_D(k+1, t)$.
- For any $t \in \mathbb{R}$ and any $k \in \mathbb{N}$, $|\lambda_D(k, t) - \lambda_{D'}(k, t)| \leq d_B(D, D')$ where $d_B(D, D')$ denotes the bottleneck distance between D and D' .

stability properties of persistence landscapes

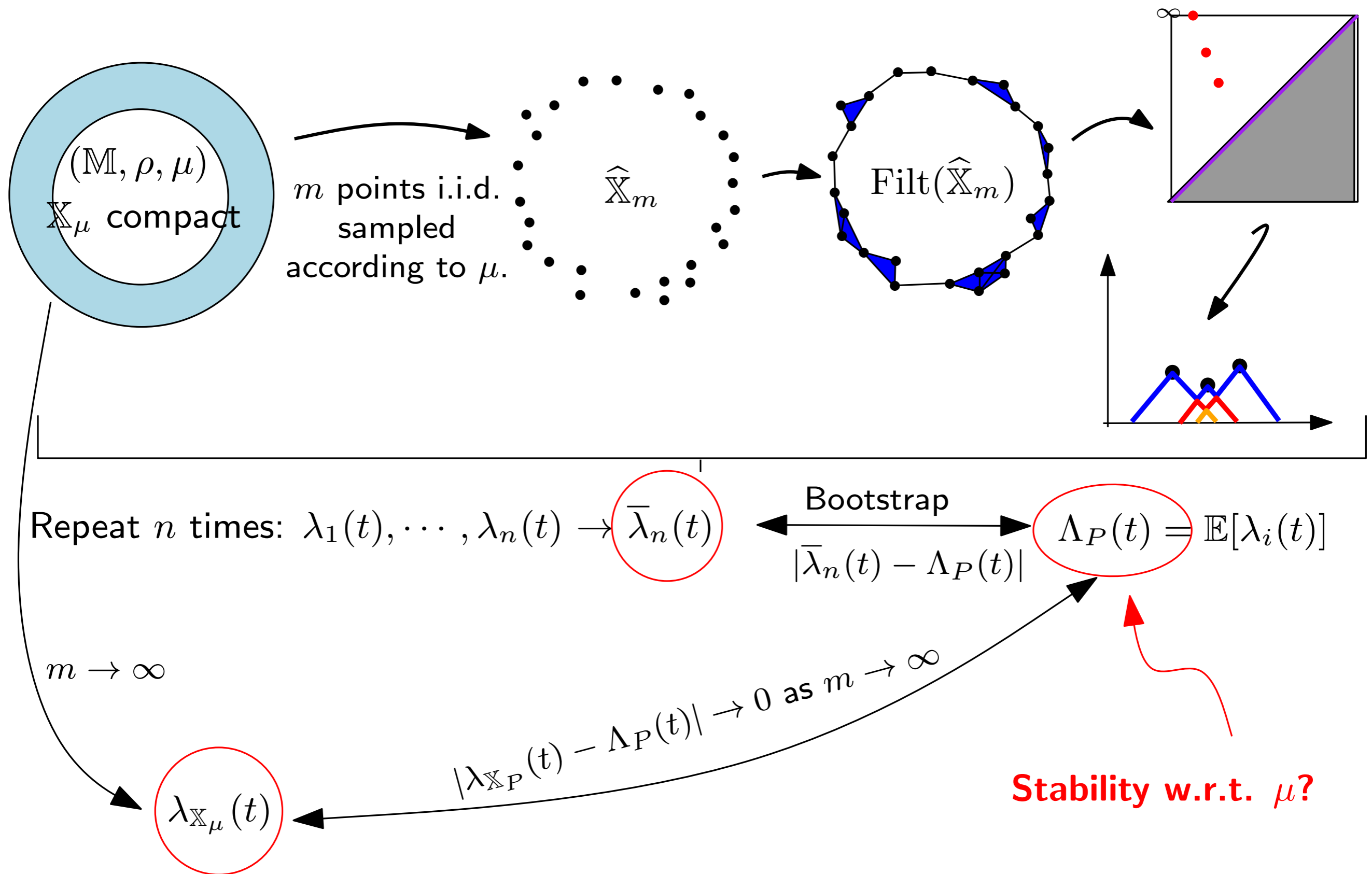
Persistence landscapes



- Persistence encoded as an element of a functional space (vector space!).
- Expectation of distribution of landscapes is well-defined and can be approximated from average of sampled landscapes.
- process point of view: convergence results and convergence rates \rightarrow confidence intervals can be computed using bootstrap.

[C., Fasy, Lecci, Rinaldo, Wasserman SoCG 2014]

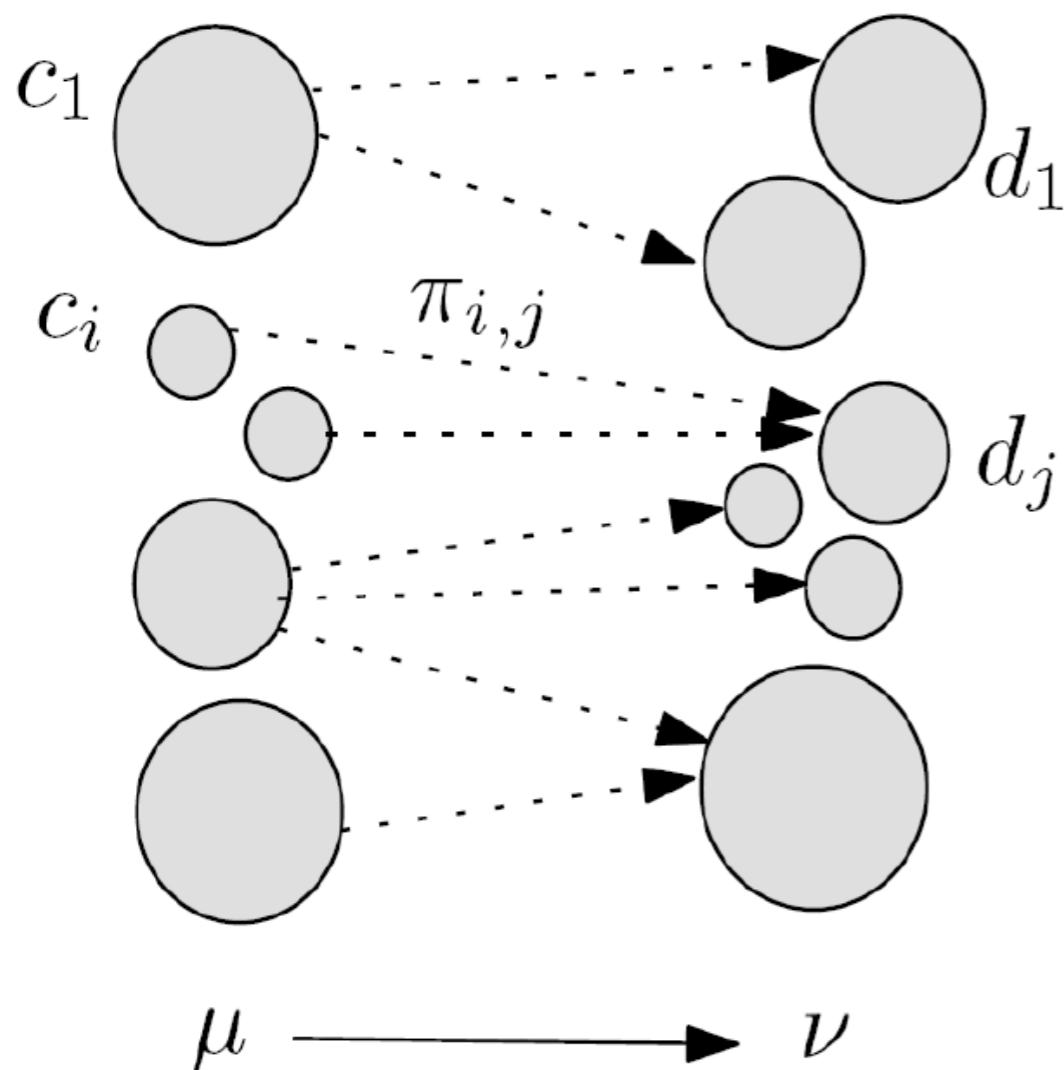
To summarize



Wasserstein distance

Let (\mathbb{M}, ρ) be a metric space and let μ, ν be probability measures on \mathbb{M} with finite p -moments ($p \geq 1$).

“The” Wasserstein distance $W_p(\mu, \nu)$ quantifies the optimal cost of pushing μ onto ν , the cost of moving a small mass dx from x to y being $\rho(x, y)^p dx$.



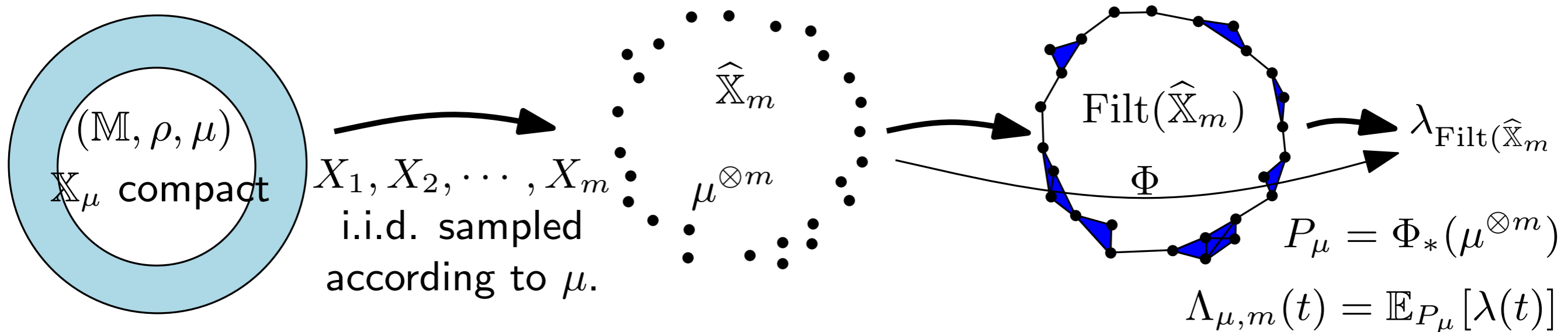
- Transport plan: Π a proba measure on $M \times M$ such that $\Pi(A \times \mathbb{R}^d) = \mu(A)$ and $\Pi(\mathbb{R}^d \times B) = \nu(B)$ for any borelian sets $A, B \subset M$.
- Cost of a transport plan:

$$C(\Pi) = \left(\int_{M \times M} \rho(x, y)^p d\Pi(x, y) \right)^{\frac{1}{p}}$$

- $W_p(\mu, \nu) = \inf_{\Pi} C(\Pi)$

(Sub)sampling and stability of expected landscapes

[C., Fasy, Lecci, Michel, Rinaldo, Wasserman ICML 2015]



Theorem: Let (\mathbb{M}, ρ) be a metric space and let μ, ν be proba measures on \mathbb{M} with compact supports. We have

$$\|\Lambda_{\mu,m} - \Lambda_{\nu,m}\|_\infty \leq m^{\frac{1}{p}} W_p(\mu, \nu)$$

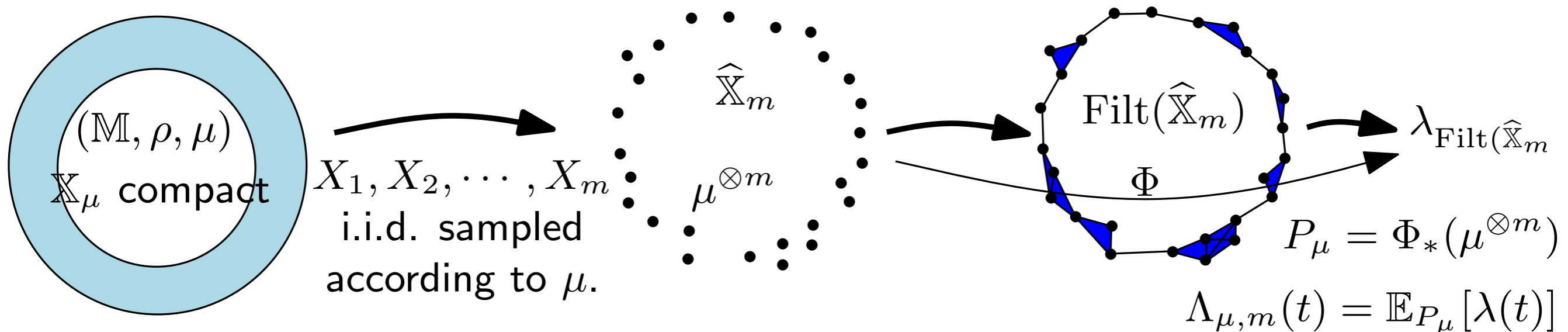
where W_p denotes the Wasserstein distance with cost function $\rho(x, y)^p$.

Remarks:

- similar results by Blumberg et al (2014) in the (Gromov-)Prokhorov metric (for distributions, not for expectations) ;
- Extended to point process setting by L. Decreusefond et al;
- $m^{\frac{1}{p}}$ cannot be replaced by a constant.

(Sub)sampling and stability of expected landscapes

[C., Fasy, Lecci, Michel, Rinaldo, Wasserman ICML 2015]



Theorem: Let (\mathbb{M}, ρ) be a metric space and let μ, ν be proba measures on \mathbb{M} with compact supports. We have

$$\|\Lambda_{\mu,m} - \Lambda_{\nu,m}\|_\infty \leq m^{\frac{1}{p}} W_p(\mu, \nu)$$

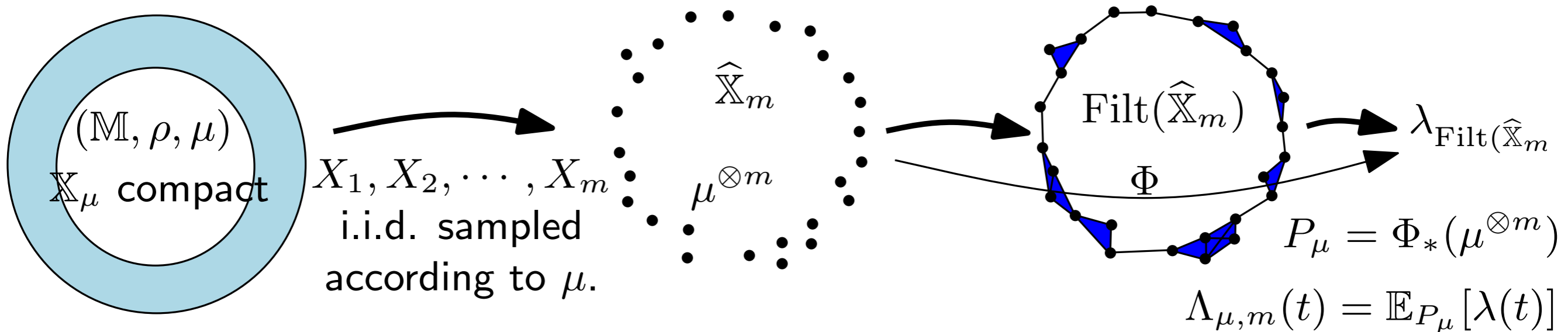
where W_p denotes the Wasserstein distance with cost function $\rho(x, y)^p$.

Consequences:

- Subsampling: efficient and easy to parallelize algorithm to infer topol. information from huge data sets.
- Robustness to outliers.
- R package TDA + Gudhi library: <https://project.inria.fr/gudhi/software/>

(Sub)sampling and stability of expected landscapes

[C., Fasy, Lecci, Michel, Rinaldo, Wasserman ICML 2015]



Theorem: Let (\mathbb{M}, ρ) be a metric space and let μ, ν be proba measures on \mathbb{M} with compact supports. We have

$$\|\Lambda_{\mu,m} - \Lambda_{\nu,m}\|_\infty \leq m^{\frac{1}{p}} W_p(\mu, \nu)$$

where W_p denotes the Wasserstein distance with cost function $\rho(x, y)^p$.

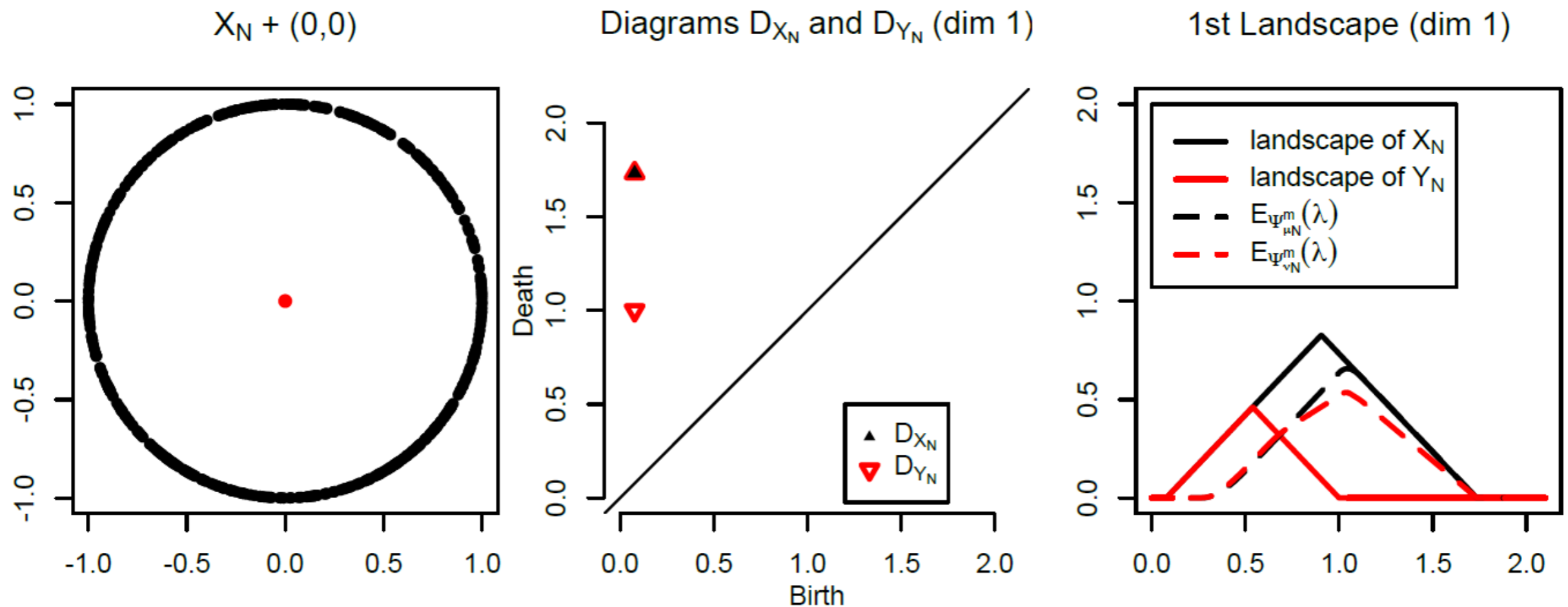
Proof:

1. $W_p(\mu^{\otimes m}, \nu^{\otimes m}) \leq m^{\frac{1}{p}} W_p(\mu, \nu)$
2. $W_p(P_\mu, P_\nu) \leq W_p(\mu^{\otimes m}, \nu^{\otimes m})$ (stability of persistence!)
3. $\|\Lambda_{\mu,m} - \Lambda_{\nu,m}\|_\infty \leq W_p(P_\mu, P_\nu)$ (Jensen's inequality)

(Sub)sampling and stability of expected landscapes

[C., Fasy, Lecci, Michel, Rinaldo, Wasserman ICML 2015]

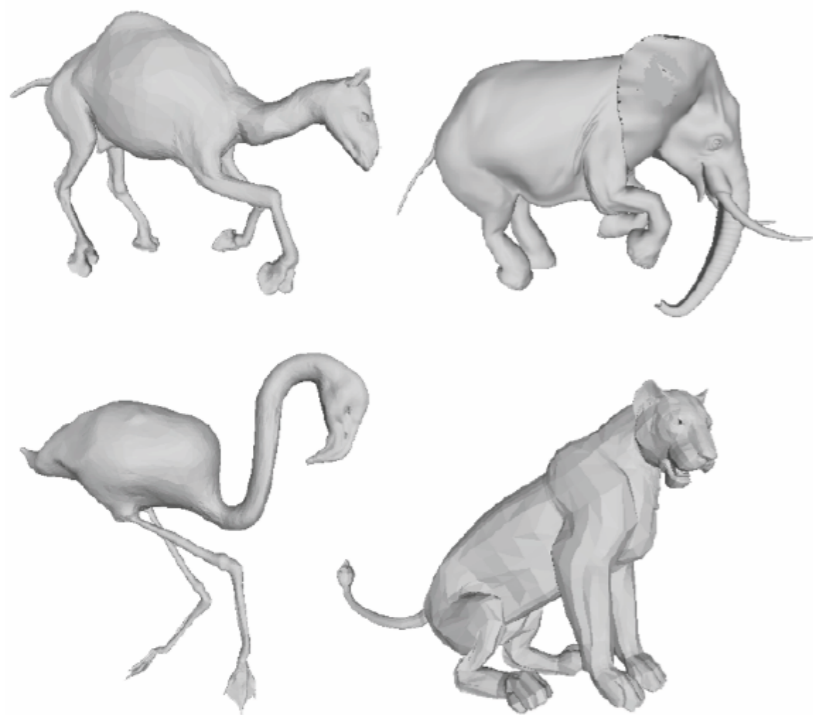
Example: Circle with one outlier.



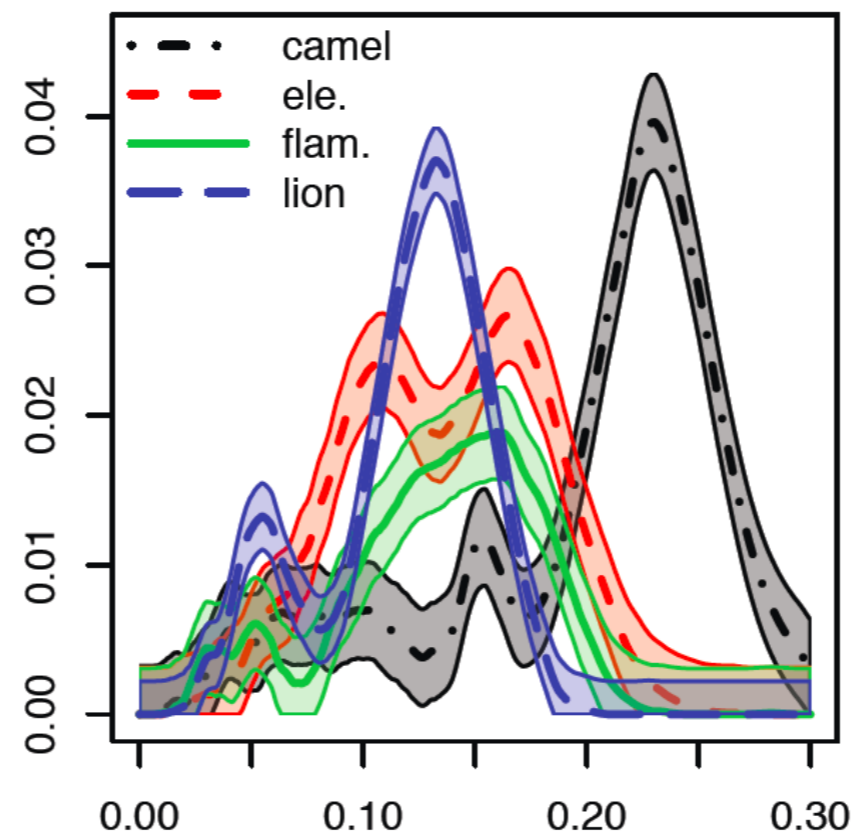
(Sub)sampling and stability of expected landscapes

[C., Fasy, Lecci, Michel, Rinaldo, Wasserman ICML 2015]

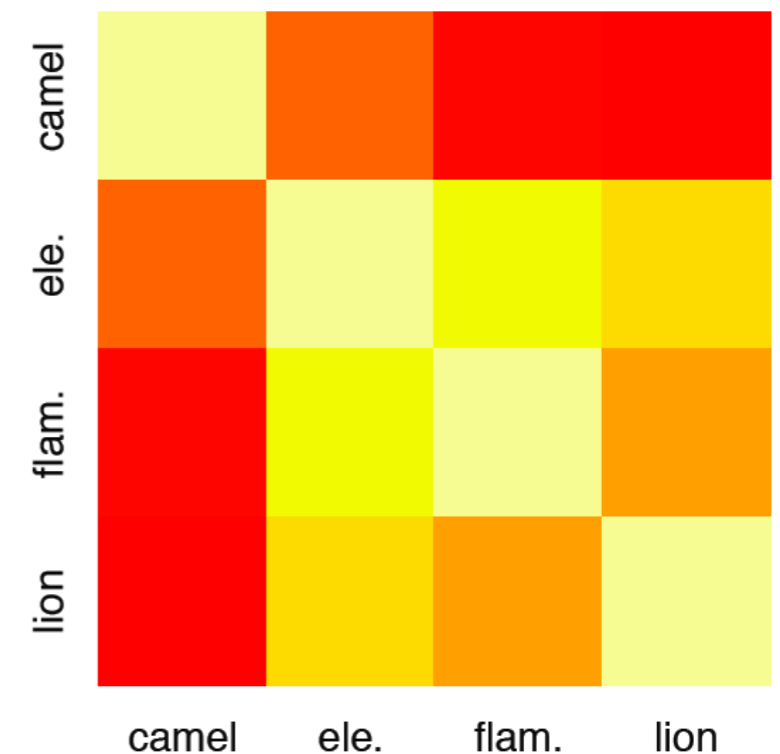
Example: 3D shapes



Average Landscapes



Dissimilarity Matrix

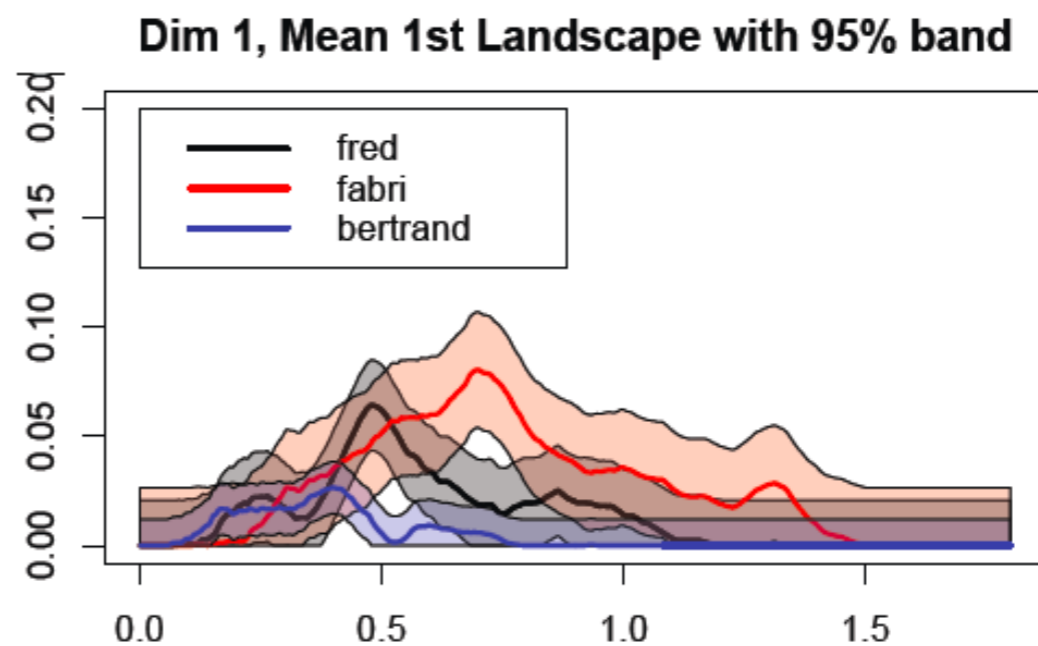
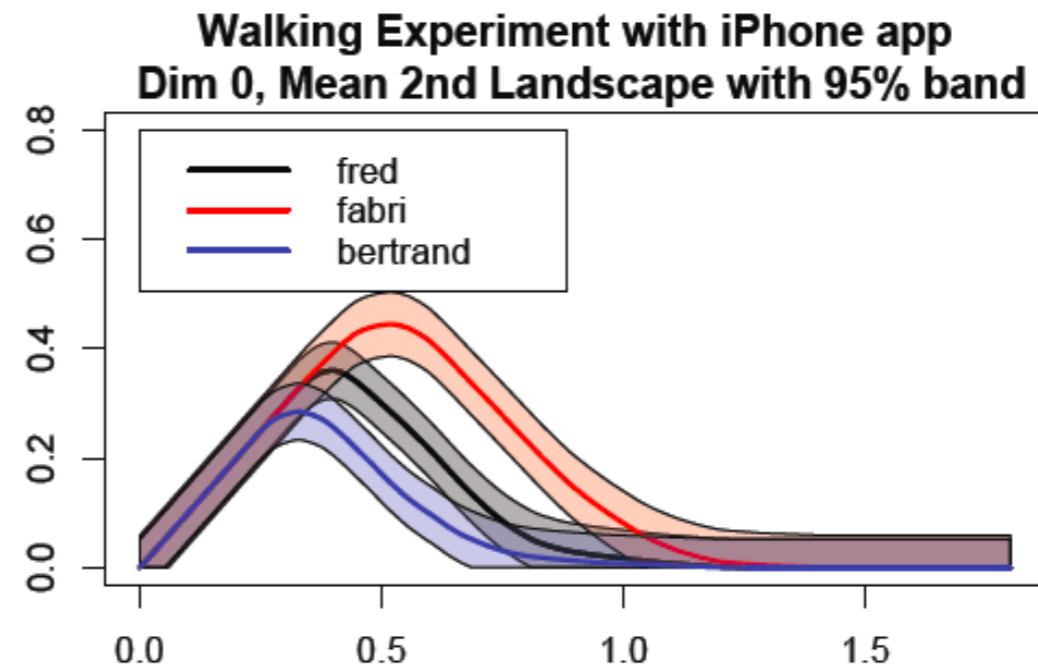
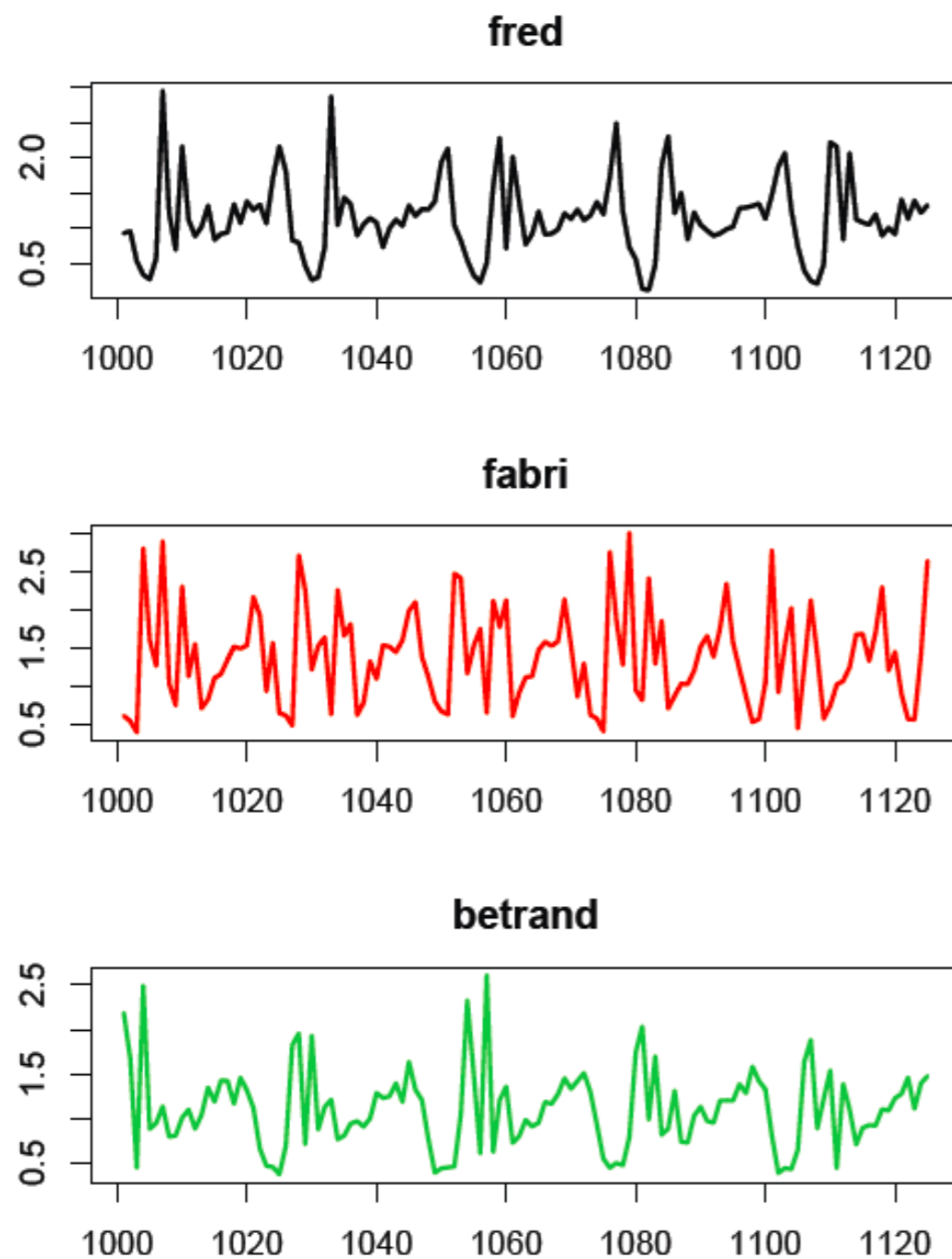


From $n = 100$ subsamples of size $m = 300$

(Sub)sampling and stability of expected landscapes

[C., Fasy, Lecci, Michel, Rinaldo, Wasserman ICML 2015]

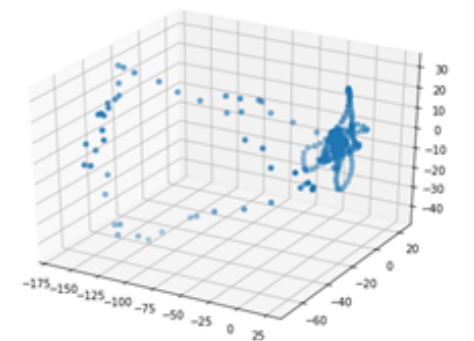
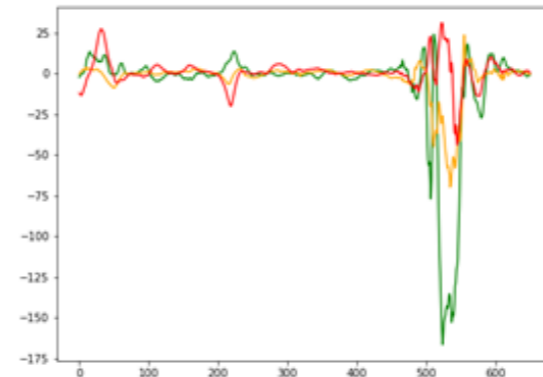
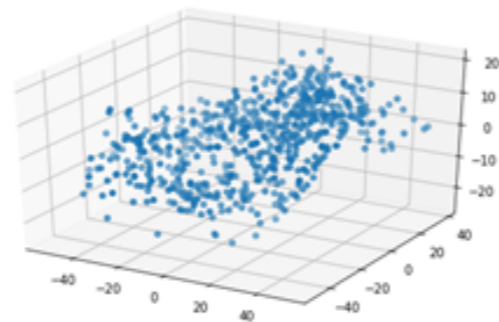
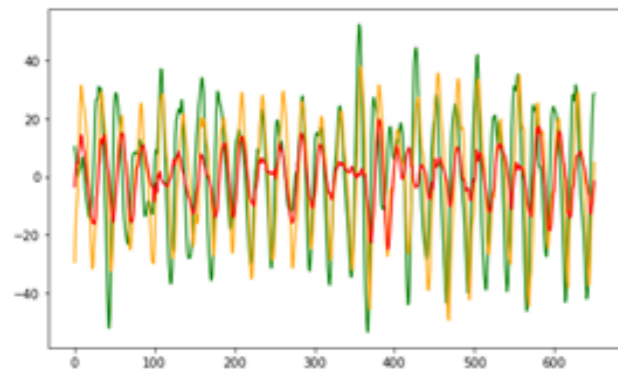
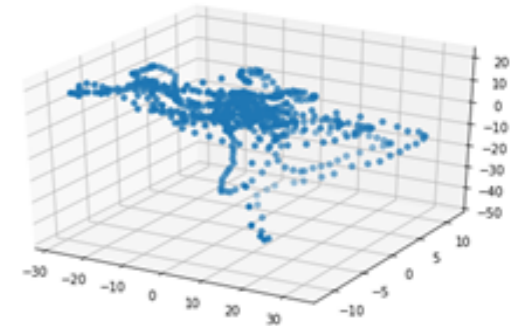
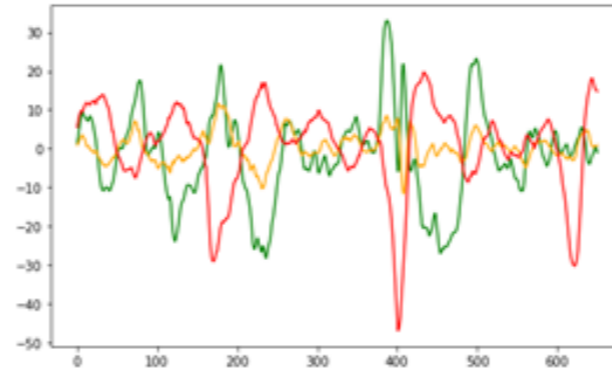
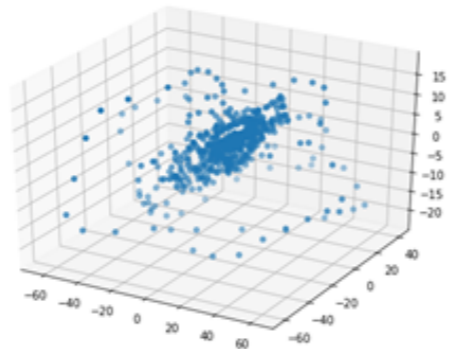
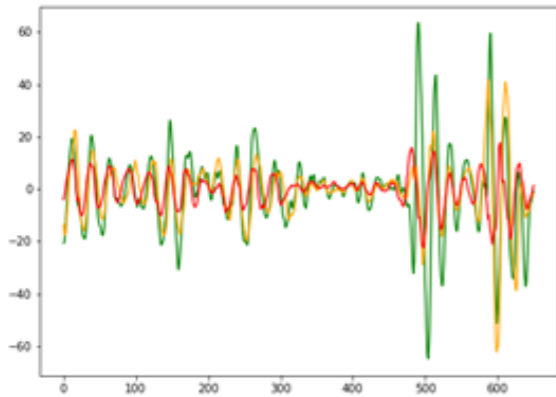
(Toy) Example: Accelerometer data from smartphone.



- spatial time series (accelerometer data from the smartphone of users).
- no registration/calibration preprocessing step needed to compare!

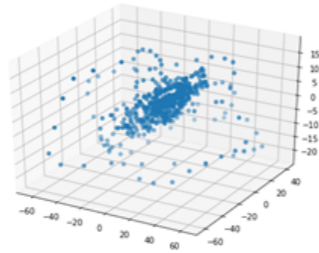
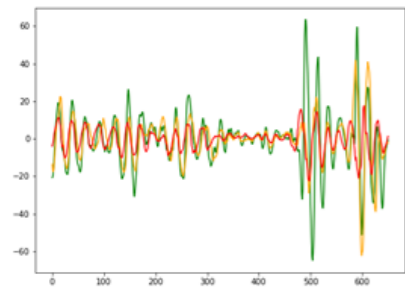
TDA and Machine Learning:
some illustrative examples on real applications

TDA and Machine Learning for sensor data

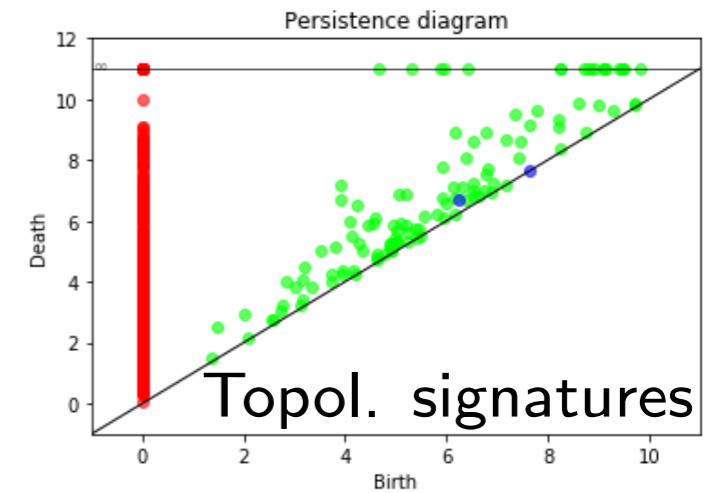


(Multivariate) time-dependent data can be converted into point clouds:
sliding window, time-delay embedding,...

TDA and Machine Learning for sensor data



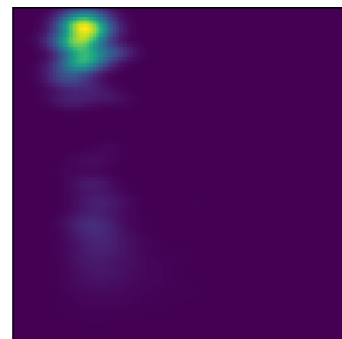
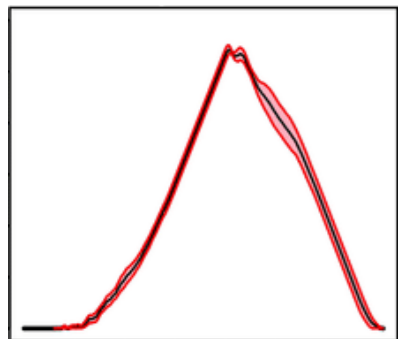
TDA pipeline
GUDHI
software



Feature engineering



Representations of persistence (linearization):



...



ML/AI

Features extraction
Random forests
Deep learning
Etc...

combined with other features!

Persistent silhouette
[Chazal & al, 2013]

Persistent surface
[Adams & al, 2016]

With landscapes: patient monitoring

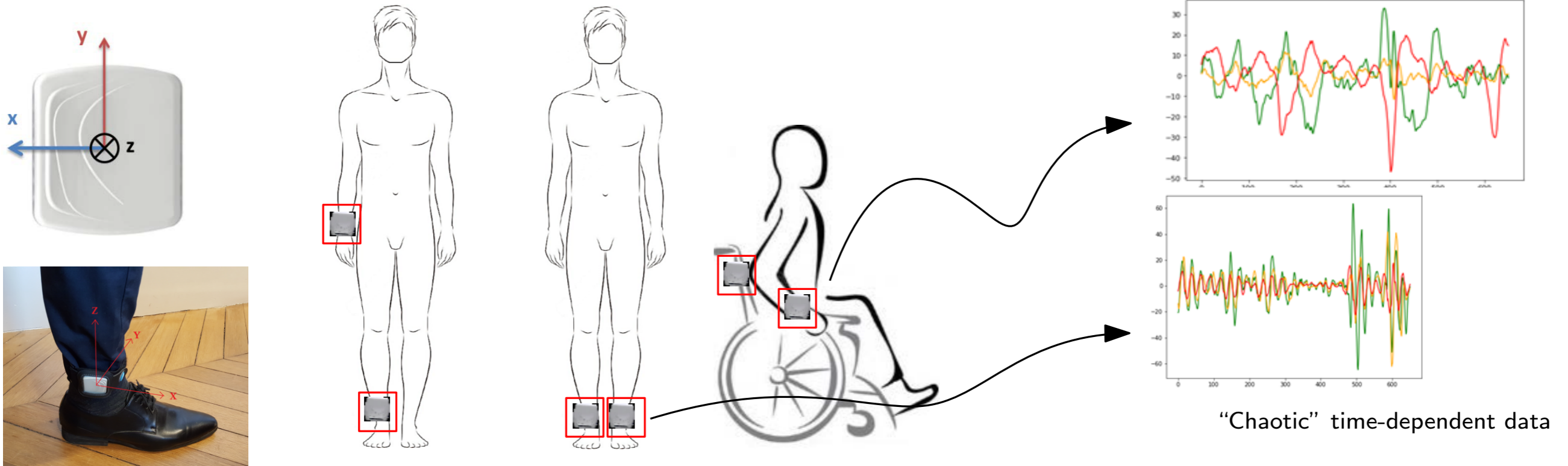
A joint industrial research project between



and



A French SME with innovating technology to reconstruct pedestrian trajectories from inertial sensors (ActiMyo)

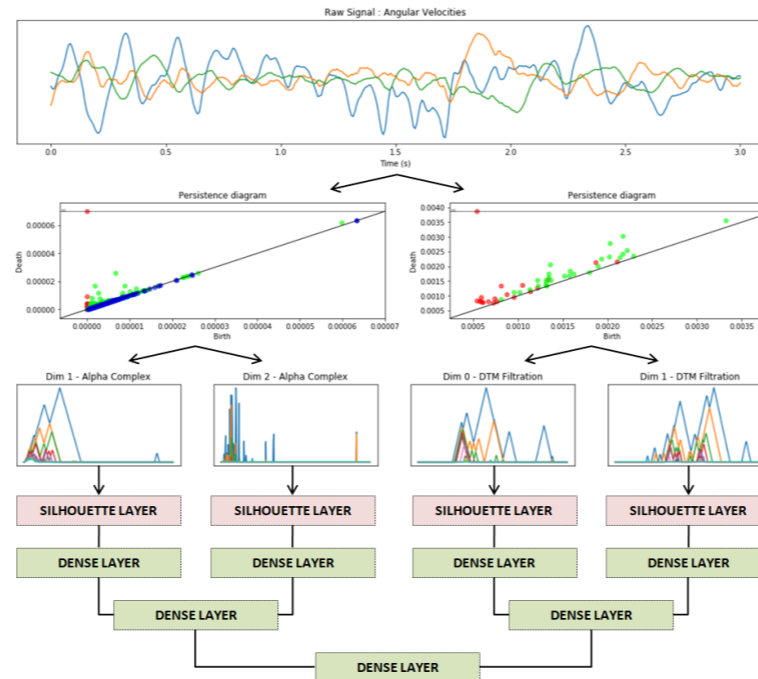
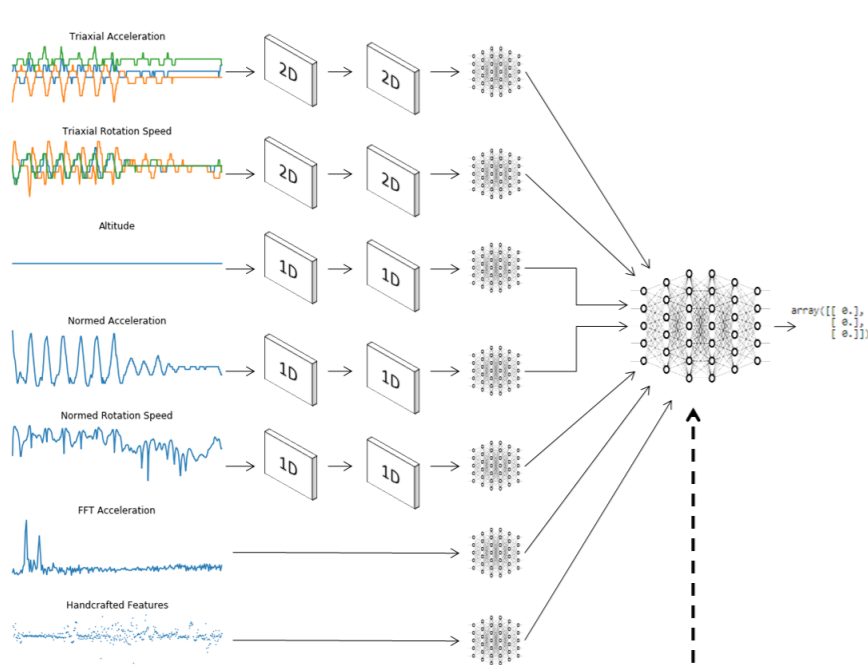


Objective: precise analysis of movements and activities of pedestrians.

Applications: personal healthcare; medical studies; defense.

With landscapes: patient monitoring

Example: Dyskinesia crisis detection and activity recognition:



Class	Naive	Multi	FEA	QUA	TDA
Walking	97.6	98.4	99.3	99.0	99.5
Upstairs	97.2	99.8	97.8	98.0	97.7
Downstairs	99.6	99.7	99.0	98.4	98.3
Sitting	87.1	93.1	89.7	91.8	96.5
Standing	87.0	97.7	97.2	97.2	98.1
Laying	92.4	100.	99.8	99.9	100.
Stand-Sit	90.8	95.6	89.1	91.3	93.4
Sit-Stand	100.	99.9	100.	100.	100.
Sit-Lie	87.1	81.1	84.2	90.0	95.1
Lie-Sit	81.4	81.8	85.9	91.8	87.9
Stand-Lie	74.2	87.6	86.5	87.4	81.5
Lie-Stand	80.4	72.1	83.2	77.7	83.2

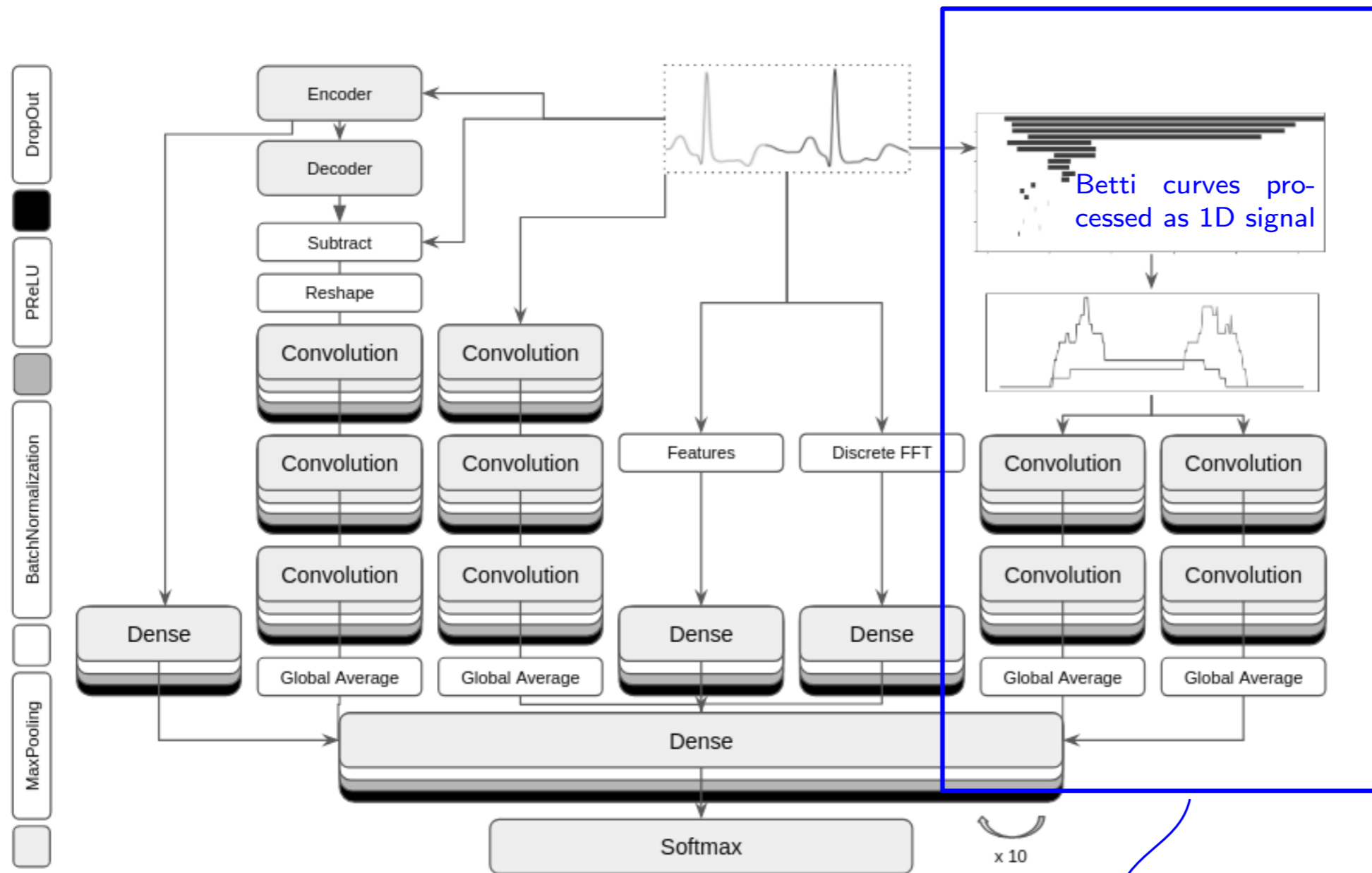
Multi-channels CNN + TDA neural network

Results on publicly available data set (HAPT) - improve the state-of-the-art.

- Data collected in non controlled environments (home) are very chaotic.
- Data registration (uncertainty in sensors orientation/position).
- Reliable and robust information is mandatory.
- Events of interest are often rare and difficult to characterize.

TDA-DL pipeline for arrhythmia detection

Objective: Arrhythmia detection from ECG data.



- Improvement over state-of-the-art.
- Better generalization.

	Accuracy[%]
UCLA (2018)	93.4
Li et al. (2016)	94.6
Inria-Fujitsu (2018)*	98.6

Thank you for your attention!