

Robust Anisotropic Power-Functions-Based Filtrations for Clustering

Claire BréchetEAU

Laboratoire de Mathématiques Jean Leray & École Centrale de Nantes, France
claire.brecheteau@ec-nantes.fr

Abstract

We consider robust power-distance functions that approximate the distance function to a compact set, from a noisy sample. We pay particular interest to robust power-distance functions that are anisotropic, in the sense that their sublevel sets are unions of ellipsoids, and not necessarily unions of balls. Using persistence homology on such power-distance functions provides robust clustering schemes. We investigate such clustering schemes and compare the different procedures on synthetic and real datasets. In particular, we enhance the good performance of the anisotropic method for some cases for which classical methods fail.

2012 ACM Subject Classification Theory of computation → Unsupervised learning and clustering

Keywords and phrases Power functions, Filtrations, Hierarchical Clustering, Ellipsoids

Digital Object Identifier 10.4230/LIPIcs.SoCG.2020.23

Related Version A full version of the paper is available at <https://hal.archives-ouvertes.fr/hal-02397100>.

Supplementary Material At <https://hal.archives-ouvertes.fr/hal-02397100>, the source code is available, as an annex file.

Acknowledgements I am extremely grateful to Samuel Tapie, for his suggestion to use tangency of ellipsoids at their first intersection point, to derive the expression of their intersection radius.

1 Introduction

Often data can be represented as a point cloud \mathbb{X} in a Euclidean space \mathbb{R}^d . Grouping data into clusters as homogeneous and well-separated as possible is the purpose of clustering. When no label is known in advance, we talk about unsupervised clustering. Topological data analysis (TDA) tools are designed to understand the shape of the data. Thereby, such tools may help to understand the shape of clusters in which to group the data. In this paper, we develop and study a TDA-based unsupervised clustering scheme. In addition, our method detects and removes points that do not really belong to any cluster; the outliers.

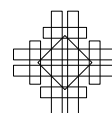
Clustering datasets is of extreme importance in multiple domains including medicine and social networks among others. The classical k -means method clusters data into isotropic clusters. In particular, the trimmed version of k -means of [14] that removes outliers, supplies balls-shaped clusters. These two algorithms have been extended by [2, 5] for Bregman-balls-shaped clusters, see also `tclust` [17] for ellipsoidal clusters. Such methods are well-suited for data generated according to mixtures of distributions which sublevel-set are Bregman balls themselves. For more general datasets, for instance, a sample of point from a disconnected manifold, these methods are no longer appropriate. Spectral clustering methods [27] perform such tasks, but are not robust to outliers. DBSCAN [19] is an algorithm based on a fixed upper-level set of an approximation of the density, and consequently, does not provide a multiscale information. Via a dendrogram, the classical single-linkage hierarchical clustering algorithm provides such a multiscale information. The dendrogram encodes information about the connectivity of unions of balls centered at points in \mathbb{X} , or equivalently, of the sublevel



© Claire BréchetEAU;
licensed under Creative Commons License CC-BY
36th International Symposium on Computational Geometry (SoCG 2020).
Editors: Sergio Cabello and Danny Z. Chen; Article No. 23; pp. 23:1–23:15



Leibniz International Proceedings in Informatics
Schloss Dagstuhl – Leibniz-Zentrum für Informatik, Dagstuhl Publishing, Germany



sets of the distance function to \mathbb{X} . For a fixed radius r , the Čech complex is a simplicial complex defined as the collection of simplices (vertex, edge, triangle, tetrahedron) for which the r -balls centered at the vertices have a non-empty common intersection. We call 1-skeleton its subcomplex (a graph) that contains only vertices and edges. The non-decreasing family of such graphs indexed by $r \in \mathbb{R}$ is called a filtration. Single-linkage is a persistence-based method since is based on the persistence, prominence or equivalently lifetime of the connected components into this graph filtration, however, it is not robust to outliers. The algorithm ToMATo in [12] is robust and persistence-based. Indeed, it is based on a graph filtration built from a neighborhood graph and a (robust) distance-like function whose values guide the appearance of vertices and edges in the graph filtration. An example of robust distance function that Chazal et al. consider in [12] is given by the distance-to-measure (DTM) [10]. Note that the graph is a priori not intrinsic to the distance function, which may cause bad clustering. For instance, an edge that links two vertices with small distance-function value but intersects an area with large distance function value, may link two clusters that should not be. This problem was the cause of failure of the single-linkage method for data corrupted by outliers. Alternative filtrations that do not suffer from this problem are the DTM-filtration [1], or the power filtrations [7], based on the 1-skeleton of the Čech filtration associated to the sublevel sets of a power distance function: a function of type $x \mapsto \min_{i \in I} \|x - m_i\|^2 + \omega_i$ for some $(m_i)_{i \in I}$ in \mathbb{R}^d and $(\omega_i)_{i \in I}$ in \mathbb{R} . Some approximations of the DTM that are power functions have been introduced and studied in the literature: the k -witnessed distance [18], the power distance [7], the c -PDTM [6] whose sublevel sets are unions of c balls, and the c -PLM [4] whose sublevel sets are unions of c ellipsoids, with c possibly much smaller than the sample size. The last two functions are robust to outliers since their construction is based on the principle of trimmed least squares [26].

Contributions

By replacing balls with ellipsoids, we enlarge the notion of weighted Čech filtration into the anisotropic weighted Čech filtration. We derive an expression for the radius of intersection of two ellipsoids. We introduce a clustering algorithm based on persistence. Such a clustering algorithm can be run from any graph filtration, in particular, from the 1-skeleton of the anisotropic weighted Čech filtration, which corresponds to the filtration of sublevel sets of an anisotropic power function. We experiment this algorithm on the filtration of the c -PLM [4].

Practical interests

A clustering algorithm based on the persistence filtration of the sublevel sets of a power function is pertinent since unlike ToMATo, the graph is intrinsic to the distance function. So, no additional parameters are required for the algorithm. The main advantage of using an anisotropic power function is that its sublevel sets are ellipsoids. Much less ellipsoids are required than balls to Hausdorff-approximate a compact manifold with intrinsic dimension smaller than the ambient dimension. The clustering scheme can also be applied to decompose a set of points generated on a polygonal line into segments. Once the ellipsoids computed, the persistence algorithm runs fast. Its complexity in terms of number of comparisons is at worst $O(c^4)$, with c , the number of ellipsoids, which is much smaller than the sample size. Most importantly, the robustness of the persistence algorithm relies on the robustness of the distance function. The c -PLM [4] is robust to outliers. The guaranty for the clustering method follows from the $\|\cdot\|_\infty$ -distance closeness between the power distance function and the distance function to the underlying manifold \mathcal{X} , relatively to the minimal distance between the connected components of \mathcal{X} . Note that such a proximity condition is sufficient but not necessary, as illustrated by the different numerical examples, with the c -PLM.

Organisation of the paper

In Section 2, we recall the notions of power function and weighted Čech filtration, the filtration of the nerves of its sublevel sets, that we extend to anisotropic power functions. We prove some stability and approximation properties for such filtrations. Examples of robust power filtrations are also displayed. The main clustering algorithm, Algorithm 1 is given in Section 3. This algorithm applies to any filtration of graphs, including the graph filtrations obtained as the 1-skeleton of a weighted Čech filtration. We enumerate other types of filtrations that fit into this framework. Finally, we implement Algorithm 1 with the robust anisotropic aforementioned power function in Section 4. We compare this method to other clustering methods on synthetic and real datasets.

2 Power-functions-based filtrations for robust clustering

In the sequel, we will recall the notion of filtration for subsets of \mathbb{R}^d (equipped with the Euclidean norm $\|\cdot\|$) and for simplicial complexes. We will consider a class of functions for which filtrations associated to sublevel sets are easily represented by filtrations of simplicial complexes, making the evolution of their connected components tractable: the power functions. In addition, we will give an example of robust power-functions [6] that can be built from a probability distribution or a pointset \mathbb{X} . Their sublevel sets are unions of c balls, with c possibly much smaller than the size of \mathbb{X} . Most importantly, we will also give an example of a robust anisotropic power-function, whose sublevel sets are unions of c ellipsoids [4]. Both of these power functions will be considered in the next sections for clustering purposes.

2.1 Generalities on filtrations

A filtration indexed by a time set $T \subset \mathbb{R}$ is a family $(V^t)_{t \in T}$ of subsets of \mathbb{R}^d , non-decreasing for the inclusion (i.e. $\forall t \leq t', V^t \subset V^{t'}$). A typical example is the filtration of the sub-level sets of a function $f : \mathbb{R}^d \mapsto \mathbb{R}$, $(f^{-1}((-\infty, t]))_{t \in T}$. For any simplex S with finite vertex set \mathbb{X} , a filtration of simplicial complexes of S is a non-decreasing family $(S^t)_{t \in T}$ of subcomplexes of S , meaning that for every $t \leq t'$, any simplex of S^t is also a simplex of $S^{t'}$.

The interleaving pseudo-distance between two filtrations $(V^t)_{t \in T}$ and $(W^t)_{t \in T}$ is defined as the smallest $\epsilon > 0$ such that $(V^t)_{t \in T}$ and $(W^t)_{t \in T}$ are ϵ -interleaved, i.e. such that: $\forall t \in T, V^t \subset W^{t+\epsilon}$ and $W^t \subset V^{t+\epsilon}$. This definition extends to simplicial complexes. Note that the sub-level-sets filtrations of two functions f and g satisfying $\|f - g\|_\infty \leq \epsilon$ are ϵ -interleaved. We will see in Section 3 that the notion of interleaving is primordial, since it measures the difference of topology between two filtrations. In particular, the stability of our sub-level-sets-based clustering scheme will be guaranteed from the closeness of the functions.

2.2 Power-functions-based filtrations

In this paper, we consider classes of functions whose sub-level sets filtration has a sparse representation, the power functions. The sublevel sets of these functions can be represented by simplicial complexes in so-called weighted Čech filtrations. We will consider two types of power functions, the isotropic and the anisotropic ones.

2.2.1 The isotropic case

An isotropic power function is a function $f_{\mathbf{m}, \boldsymbol{\omega}} : \mathbb{R}^d \rightarrow \mathbb{R}$ defined from an index set $I = \llbracket 1, c \rrbracket$, a family of centers $\mathbf{m} = (m_i)_{i \in I}$ in \mathbb{R}^d and a family of weights $\boldsymbol{\omega} = (\omega_i)_{i \in I}$ in \mathbb{R} by $f_{\mathbf{m}, \boldsymbol{\omega}} : x \mapsto \min_{i \in I} \|x - m_i\|^2 + \omega_i$. A simple example of power function is the squared

Euclidean distance function to a set of points \mathbb{X} , $d_{\mathbb{X}}^2 : x \in \mathbb{R}^d \mapsto \min_{m \in \mathbb{X}} \|x - m\|^2$. The sublevel sets of $f_{\mathbf{m}, \boldsymbol{\omega}}$, $V_{\mathbf{m}, \boldsymbol{\omega}}^t = f_{\mathbf{m}, \boldsymbol{\omega}}^{-1}((-\infty, t])$, are unions of at most c balls $\mathcal{B}_i^t = \overline{B}(m_i, \sqrt{t - \omega_i})$ with $\overline{B}(m, r) = \{x \in \mathbb{R}^d \mid \|x - m\| \leq r\}$. Note that \mathcal{B}_i^t is empty for $t < \omega_i$ and two balls \mathcal{B}_i^t and \mathcal{B}_j^t intersect if and only if $t \geq t_{i,j}$ with $t_{i,j} = \frac{(\omega_j - \omega_i)^2 + 2(\omega_j + \omega_i)\|m_j - m_i\|^2 + \|m_j - m_i\|^4}{4\|m_j - m_i\|^2}$. The connectivity of $V_{\mathbf{m}, \boldsymbol{\omega}}^t$ can be encoded in a graph $\mathcal{G}_{\mathbf{m}, \boldsymbol{\omega}}^t$, whose vertices are indices $i \in I$ such that $\omega_i \leq t$ and whose edges are pairs of vertices $[i, j]$ such that $t_{i,j} \leq t$. Indeed, $\mathcal{G}_{\mathbf{m}, \boldsymbol{\omega}}^t$ and $V_{\mathbf{m}, \boldsymbol{\omega}}^t$ have the same number of connected components, and m_i and m_j are in the same connected component in $V_{\mathbf{m}, \boldsymbol{\omega}}^t$ if and only if i and j are also in the same component in $\mathcal{G}_{\mathbf{m}, \boldsymbol{\omega}}^t$.

More generally, the topological information of $V_{\mathbf{m}, \boldsymbol{\omega}}^t$ (number of connected components, loops, voids etc.) can be encoded in the weighted Čech complex $\text{Cech}_{\mathbf{m}, \boldsymbol{\omega}}(t)$, defined as the nerve of the union of balls $(\mathcal{B}_i^t)_{i \in I}$: $\text{Cech}_{\mathbf{m}, \boldsymbol{\omega}}(t) = \{\sigma \subset I \mid \bigcap_{i \in \sigma} \mathcal{B}_i^t \neq \emptyset\}$, [1, 7, 3]. According to the Nerve Lemma [20, Corollary 4G.3], any sublevel set $V_{\mathbf{m}, \boldsymbol{\omega}}^t$ is homotopic to $\text{Cech}_{\mathbf{m}, \boldsymbol{\omega}}(t)$ and thus contains the same topological information. For computational reasons, the weighted Vietoris-Rips filtration is frequently considered as a provably good surrogate for the weighted Čech filtration $(\text{Cech}_{\mathbf{m}, \boldsymbol{\omega}}(t))_{t \in T}$. The weighted Vietoris-Rips complex $\text{VR}_{\mathbf{m}, \boldsymbol{\omega}}(t)$ is the flag complex of $\mathcal{G}_{\mathbf{m}, \boldsymbol{\omega}}^t$ ($\mathcal{G}_{\mathbf{m}, \boldsymbol{\omega}}^t$ is the 1-skeleton of the weighted Čech complex): $\text{VR}_{\mathbf{m}, \boldsymbol{\omega}}(t) = \{\sigma \subset I \mid \forall i, j \in \sigma, \mathcal{B}_i^t \cap \mathcal{B}_j^t \neq \emptyset\}$. Indeed, as a direct consequence of [3, Theorem 3.2] which is a generalization of the non-weighted case in [15, Theorem 2.5.], if the weights in $\boldsymbol{\omega}$ are non-negative, then these two filtrations are interleaved:

$$\forall 0 < t' \leq \frac{d+1}{2d}t, \text{VR}_{\mathbf{m}, \boldsymbol{\omega}}(t') \subset \text{Cech}_{\mathbf{m}, \boldsymbol{\omega}}(t) \subset \text{VR}_{\mathbf{m}, \boldsymbol{\omega}}(t). \quad (1)$$

These notions can all be extended to anisotropic power functions.

2.2.2 The anisotropic case

Consider $I = \llbracket 1, c \rrbracket$, centers $\mathbf{m} = (m_i)_{i \in I}$ in \mathbb{R}^d , weights $\boldsymbol{\omega} = (\omega_i)_{i \in I}$ in \mathbb{R} and matrices $\boldsymbol{\Sigma} = (\Sigma_i)_{i \in I}$ in \mathcal{M}_d , the set of definite positive symmetric matrices. An anisotropic power function is a function $f_{\mathbf{m}, \boldsymbol{\omega}, \boldsymbol{\Sigma}} : \mathbb{R}^d \rightarrow \mathbb{R}$ defined from I , \mathbf{m} , $\boldsymbol{\omega}$ and $\boldsymbol{\Sigma}$ by $f_{\mathbf{m}, \boldsymbol{\omega}, \boldsymbol{\Sigma}} : x \mapsto \min_{i \in I} \|x - m_i\|_{\Sigma_i}^2 + \omega_i$. For any matrix $\Sigma \in \mathcal{M}_d$ and $x \in \mathbb{R}^d$, $\|x\|_{\Sigma^{-1}} = \sqrt{x^T \Sigma^{-1} x}$ denotes the Σ -Mahalanobis norm of x . The sublevel sets of $f_{\mathbf{m}, \boldsymbol{\omega}, \boldsymbol{\Sigma}}$, $V_{\mathbf{m}, \boldsymbol{\omega}, \boldsymbol{\Sigma}}^t = f_{\mathbf{m}, \boldsymbol{\omega}, \boldsymbol{\Sigma}}^{-1}((-\infty, t])$, are unions of at most c ellipsoids $\mathcal{E}_i^t = \overline{B}_{\Sigma_i}(m_i, \sqrt{t - \omega_i}) = \{x \in \mathbb{R}^d \mid \|x - m_i\|_{\Sigma_i}^2 \leq t - \omega_i\}$. Again, \mathcal{E}_i^t is empty for $t < \omega_i$ and the intersection time $t_{i,j}$ of \mathcal{E}_i^t and \mathcal{E}_j^t is given below. The relative question of the emptiness of the intersection of two ellipsoids is tackled in [28, 25].

► **Proposition 1.** *Consider two ellipsoids $\mathcal{E}_i^t = \overline{B}_{\Sigma_i}(m_i, \sqrt{t - \omega_i})$ and $\mathcal{E}_j^t = \overline{B}_{\Sigma_j}(m_j, \sqrt{t - \omega_j})$ with $\omega_i \leq \omega_j$ in \mathbb{R} , m_i and m_j in \mathbb{R}^d , $\Sigma_i = P_i D_i P_i^T$ and $\Sigma_j = P_j D_j P_j^T$ in \mathcal{M}_d , with two positive diagonal matrices D_i and D_j and two orthogonal matrices P_i and P_j from the spectral theorem. Set $\tilde{\Sigma} = \sqrt{D_i} P_i^T \Sigma_j^{-1} P_i \sqrt{D_i} = \tilde{P} \tilde{D} \tilde{P}^T$, for orthogonal and diagonal matrices \tilde{P} and $\tilde{D} = \text{diag}(\lambda_1, \lambda_2, \dots, \lambda_d)$, and $\tilde{m} = \tilde{P}^T \sqrt{D_i^{-1} P_i^T} (m_j - m_i)$. Ellipsoids \mathcal{E}_i^t and \mathcal{E}_j^t intersect if and only if $t \geq t_{i,j}$ for $t_{i,j} = \omega_j$ when $\|\tilde{m}\| \leq \sqrt{\omega_j - \omega_i}$, and $t_{i,j} = \omega_j + \sum_{k=1}^d \left(\frac{\lambda_k \tilde{m}_k}{\lambda + \lambda_k} \right)^2 \lambda_k$ when $\|\tilde{m}\| > \sqrt{\omega_j - \omega_i}$. The positive number λ is the unique solution of the following equation:*

$$\sum_{k=1}^d \frac{\lambda_k - \lambda^2}{(\lambda + \lambda_k)^2} \lambda_k \tilde{m}_k^2 = \omega_j - \omega_i. \quad (2)$$

The proof is based on the fact that the ellipsoids \mathcal{E}_i^t and \mathcal{E}_j^t are tangent at their first intersection point, and the corresponding gradients are collinear. In the context of isotropy (i.e. for $\Sigma_i = \Sigma_j = I_d$, the identity matrix of \mathbb{R}^d) $\tilde{m} = m_j - m_i$, and when $\|m_j - m_i\| > \sqrt{\omega_j - \omega_i}$, (2) has a unique positive solution given by $\lambda = \frac{\omega_i - \omega_j + \|m_j - m_i\|^2}{\omega_j - \omega_i + \|m_j - m_i\|^2}$. We recover the merging time $t_{i,j}$ given in Section 2.2.1. Now, define $\mathcal{G}_{\mathbf{m},\omega,\Sigma}^t$, $\text{Cech}_{\mathbf{m},\omega,\Sigma}(t)$ and $\text{VR}_{\mathbf{m},\omega,\Sigma}(t)$, the anisotropic counterparts of $\mathcal{G}_{\mathbf{m},\omega}^t$, $\text{Cech}_{\mathbf{m},\omega}(t)$ and $\text{VR}_{\mathbf{m},\omega}(t)$. The nerve lemma still applies, since unions of ellipsoids are contractible. Although this paper is mostly based on the study of connected components for clustering, anisotropic weighted Čech and Vietoris-Rips filtrations are primordial to have a tractable estimation of the topology of compact sets from suitable approximations as finite unions of ellipsoids. In fact, as their isotropic counterparts (1), these filtrations are interleaved, provided that the eigenvalues of the matrices in Σ are positive.

► **Proposition 2.** *If ω is a set on non-negative weights in \mathbb{R} and Σ a family of matrices with eigenvalues in $[\lambda_{\min}, \lambda_{\max}]$ for some $\lambda_{\min} > 0$, then for every $t > 0$ and $0 < t' \leq \frac{\lambda_{\min}}{\lambda_{\max}} \frac{d+1}{2d} t$,*

$$\text{VR}_{\mathbf{m},\omega,\Sigma}(t') \subset \text{Cech}_{\mathbf{m},\omega,\Sigma}(t) \subset \text{VR}_{\mathbf{m},\omega,\Sigma}(t). \tag{3}$$

The condition of non-negative weights is not too restrictive since for general weights, it suffices to replace ω , t and t' by $\omega - \min_{i \in I} \omega_i$, $t - \min_{i \in I} \omega_i$ and $t' - \min_{i \in I} \omega_i$ in the proposition. Then, the condition on t' becomes $\min_{i \in I} \omega_i < t' \leq \frac{\lambda_{\min}}{\lambda_{\max}} \frac{d+1}{2d} t + \left(1 - \frac{\lambda_{\min}}{\lambda_{\max}} \frac{d+1}{2d}\right) \min_{i \in I} \omega_i$. As noted in [15], when λ_{\min} equals λ_{\max} and the weights in ω are null, the term $\frac{\lambda_{\min}}{\lambda_{\max}} \frac{d+1}{2d}$ is optimal. When \mathbf{m} is the set of vertices of a regular d -simplex, the left inclusion is an equality.

Often, less ellipsoids than balls are required to describe a compact set \mathcal{X} , for a fixed level of precision (e.g. for the Hausdorff distance). For instance, a segment in \mathbb{R}^2 , and more generally, any d' -dimensional submanifold in \mathbb{R}^d , with $d' < d$. For this reason, anisotropic Čech and Vietoris-Rips filtrations are pertinent tools to compute and store the topological information about \mathcal{X} efficiently. The requisite condition is that we dispose of an anisotropic power function that is a good approximation of $d_{\mathcal{X}}^2$. Such examples of functions follow.

2.3 Examples of filtrations based on robust power functions

2.3.1 Isotropic robust power functions

Set \mathbb{X} , a set of n points generated on the neighborhood of a compact subset \mathcal{X} of \mathbb{R}^d . In order to face the non robustness of the distance function to \mathbb{X} , $d_{\mathbb{X}}$, Chazal et al. have introduced the notion of distance-to-measure (DTM), in [10]. The DTM is a counterpart of $d_{\mathbb{X}}$ robust to noise and outliers. Its robustness follows from some parameter $k \in \llbracket 1, n \rrbracket$, the number of nearest-neighbors X^1, X^2, \dots, X^k of x in \mathbb{X} , used to estimate $d_{\mathbb{X}}(x)$. The DTM $d_{\mathbb{X},k}$ is defined by $d_{\mathbb{X},k}^2 : x \mapsto \frac{1}{k} \sum_{i=1}^k \|x - X^i\|^2 = \|x - m_{x,k}\|^2 + v_{x,k}$ with $m_{x,k} = \frac{1}{k} \sum_{i=1}^k X^i$, the mean of the k nearest neighbours of x in \mathbb{X} and $v_{x,k} = \frac{1}{k} \sum_{i=1}^k \|X^i - m_{x,k}\|^2$ their variance. Note that $d_{\mathbb{X},1}$ coincides with $d_{\mathbb{X}}$ and is not robust, whereas $d_{\mathbb{X},n}(x)$ is the distance of x to the barycenter of the point cloud \mathbb{X} , up to some factor, which is robust, but very poor in terms of topological information. The DTM is actually a weighted power function [18]:

$$d_{\mathbb{X},k}^2(x) = \inf_{y \in \mathbb{R}^d} \|x - m_{y,k}\|^2 + v_{y,k}. \tag{4}$$

This follows from the fact that the mean distance between x and its k nearest neighbors is not larger than the mean distance between x and the k nearest neighbors of any other point $y \in \mathbb{R}^d$. This infimum is actually a minimum over a set of c points $\mathbf{y} = (y_i)_{i \in \llbracket 1, c \rrbracket}$ in \mathbb{R}^d , with

c of order $\binom{n}{k}$. A power approximation of the DTM, the k -witnessed distance, was defined in [18] by replacing \mathbb{R}^d by \mathbb{X} in (4). Its sublevel sets are unions of n balls. An approximation of the DTM with c (possibly much smaller than n) balls, the c -PDTM, was defined in [6], by replacing \mathbb{R}^d by a set $\mathbf{y}_{c,k}$ of c points in \mathbb{R}^d . This set $\mathbf{y}_{c,k}$ is a minimum of a “k-means”-type criterion [24], $\mathbf{y} \mapsto \sum_{i=1}^n \min_{y \in \mathbf{y}} \|X_i - m_{y,k}\|^2 + v_{y,k}$, for \mathbf{y} with cardinality c . Morally, $\mathbf{y}_{c,k}$ is chosen such that on average on \mathbb{X} , $x \mapsto \min_{y \in \mathbf{y}} \|x - m_{y,k}\|^2 + v_{y,k}$ is small. Note that the graph of the c -PDTM is necessarily above the graph of the DTM. According to [6], for a sample on a regular d' -dimensional manifold, c can be chosen of order $n^{\frac{d'}{d'+4}}$, which is much smaller than n . Moreover, the c -PDTM is a good approximation of $d_{\mathcal{X}}^2$, despite noise.

2.3.2 An anisotropic robust power function

An anisotropic version of the c -PDTM has been introduced in [4], the c -power likelihood to measure (c -PLM). It consists in replacing Euclidean norms with Mahalanobis norms. For every $x \in \mathbb{R}^d$ and $\Sigma \in \mathcal{M}_d$, set X^1, X^2, \dots, X^k the k -nearest neighbors of x in \mathbb{X} , for the Σ^{-1} -Mahalanobis norm: $\|X^i - x\|_{\Sigma^{-1}} \leq \|X^j - x\|_{\Sigma^{-1}}$ for every $i \leq j$. Denote by $m_{x,\Sigma,k}$ their mean, and by $v_{x,\Sigma,k} = \frac{1}{k} \sum_{i=1}^k \|X^i - m_{x,\Sigma,k}\|_{\Sigma^{-1}}^2$ their variance, relative to the Σ -Mahalanobis norm. Set $\theta_{c,k}$, a family of c pairs $(y, \Sigma) \in \mathbb{R}^d \times \mathcal{M}_d$ that minimizes (or which criterion is as close as possible to the optimal criterion, in case of non existence of a minimum) the following “k-means”-type criterion $R_{c,k}$ among all θ s of cardinality c : $R_{c,k}(\theta) = \sum_{i=1}^n \min_{(y,\Sigma) \in \theta} \|X_i - m_{y,\Sigma,k}\|_{\Sigma^{-1}}^2 + v_{y,\Sigma,k} + \log(\det(\Sigma))$. The term $\log(\det(\Sigma))$ prevents optimal covariance matrices to be degenerated, with Σ^{-1} going to 0. In some sense, minimizing such a criterion boils down to fit Gaussian distributions to the data set \mathbb{X} , at best. The c -PLM is the power function defined from $\theta_{c,k}$ by: $x \mapsto \min_{(y,\Sigma) \in \theta_{c,k}} \|x - m_{y,\Sigma,k}\|_{\Sigma^{-1}}^2 + v_{y,\Sigma,k} + \log(\det(\Sigma))$. A modification of the criterion $R_{c,k}$ has been introduced in [4], to remove some datapoints ($|\mathbb{X}| - sig$ for some parameter sig), when \mathbb{X} is corrupted with outliers. The criterion is given by $R_{c,k,sig}(\theta) = \min_{(i_1, i_2, \dots, i_{sig}) \in \llbracket 1, |\mathbb{X}| \rrbracket} \sum_{j=1}^{sig} \min_{(y,\Sigma) \in \theta} \|X_{i_j} - m_{y,\Sigma,k}\|_{\Sigma^{-1}}^2 + v_{y,\Sigma,k} + \log(\det(\Sigma))$.

Iterative Lloyd-type algorithms [22] provide local minima $\tilde{\theta}_{c,k}$ and $\tilde{\theta}_{c,k,sig}$ for the criteria $R_{c,k}$ and $R_{c,k,sig}$ [4]. These algorithms run in $O(ncd^2 + nkd^2 + n \log(n)c)it$ operations, with it the number of iterations of the algorithm. They consist, given $\theta = (\mathbf{y}, \Sigma)$, in splitting the space \mathbb{R}^d into weighted Σ -curved Voronoi cells, replacing centers \mathbf{y} by the centroid of the cells, and updating the matrices in Σ by a close formula from the points in the cells and ellipsoids. To compute $\tilde{\theta}_{c,k,sig}$, a trimming step is added at each iteration. For clustering, disposing of a local minimum is enough, as enhanced in the numerical illustration section, since we can remove bad centers in $\tilde{\theta}_{c,k}$ or in $\tilde{\theta}_{c,k,sig}$ with the parameter *Threshold* in Algorithm 1.

3 Persistence-based clustering from power-functions-based filtrations

3.1 Persistence for power-functions-based filtrations

Set $f_{\mathbf{m},\omega,\Sigma} : x \in \mathbb{R}^d \mapsto \min_{i \in I} \|x - m_i\|_{\Sigma^{-1}}^2 + \omega_i$, an anisotropic power-function indexed by a set $I = \llbracket 1, c \rrbracket$ and with the ω_i s sorted in non-decreasing order. As above-mentioned, the sublevel sets $V^t = f_{\mathbf{m},\omega,\Sigma}^{-1}((-\infty, t])$ are unions of at most c ellipsoids $\mathcal{E}_i^t = B_{\Sigma_i}(m_i, \sqrt{t - \omega_i})$, non empty as soon as $t \geq \omega_i$. In particular, each sublevel set of $f_{\mathbf{m},\omega,\Sigma}$ contains at most c connected components. Each connected component of V^t , V_i^t is indexed by the smallest index $i \in I$ such that m_i belongs to the component. With a language abuse, we call connected component V_i , the family of connected components $(V_i^t)_{t \in T}$ that gets born at time $t = b_i = \omega_i$ and dies at a time $t = d_i$ when V_i^t merges with another connected component V_j^t for some

$j \leq i$. Note that $d_1 = \infty$. The lifetime of the component V_i^t , $d_i - b_i$, is called persistence or prominence of the component i . This merging information is encoded in a barcode or a dendrogram. In these two representations, each line is associated to a component V_i , has length $d_i - b_i$, and begins at the height b_i . The dendrogram is obtained from the barcode by linking the bars associated to merging components, at a height given by the merging time.

When \mathbf{m} is a point set \mathbb{X} , $\Sigma_i = I_d$ and $\omega_i = 0$ for every i , clustering points accordingly to the connected components of V^t boils down to the classical single-linkage clustering procedure, with $t > 0$, calibrated in accordance with the dendrogram. This procedure is not robust to outliers. In this paper, we consider an adjacent procedure, similar to the ToMATo algorithm [12], based on the prominence of components. To be precise, in the clustering scheme, a component V_i cannot merge with another component V_j at a time t larger than $\omega_i + Stop$, for some parameter $Stop$. In other words, components with large prominence will never die in this clustering procedure. This is the purpose of Algorithm 1 in the next section.

In order to better visualize the prominence of the components, we represent their lifetimes in a persistence diagram. A persistence diagram is a multiset of points $(b_i, d_i) \in \mathbb{R}^2$ that lie above the diagonal $b = d$. Each point (b_i, d_i) is associated to a connected component V_i . The notion of persistence diagram was introduced by Edelsbrunner et al. in [16], in the broader framework of homology, and allows to compute lifetimes of additional features such as loops, voids etc. It is defined for filtrations that are regular enough, on triangulable spaces such as \mathbb{R}^d . The proper notion of regularity is the notion of q -tameness [11]. In [7, Proposition 3.5], Buchet et al. proved that the DTM is q -tame. The proof of [7] can be straightforwardly adjusted for distance functions to compact sets and most importantly, for anisotropic power functions, provided that the eigenvalues of the matrices Σ_i are all positive.

Since distance to compact sets, distance-to-measure and anisotropic power functions are q -tame, the persistence diagrams associated to their filtrations are well defined. They can be compared through the bottleneck distance, a distance between two diagrams D and D' defined as the minimal value of $\max_{x \in D, y \in D'} |y - \phi(x)|_\infty$ among functions ϕ that pair points in D with points in D' , with some points potentially paired to diagonal points. Diagrams associated to interleaved filtrations are close, according to the following theorem.

► **Theorem 3** (Stability of persistence diagrams [11, 9, 13]). *If two filtrations V and W are q -tame and ϵ -interleaved, then the persistence diagrams of these filtrations are ϵ -close in bottleneck distance.*

According to Theorem 3, the persistence diagram of any anisotropic power function $f_{\mathbf{m}, \omega, \Sigma}$ that is $\epsilon - \|\cdot\|_\infty$ close to $d_{\mathcal{X}}$ is ϵ -bottleneck close to the persistence diagram of the sublevel sets of $d_{\mathcal{X}}$. Consequently, prominence of the connected components of \mathcal{X} can be deduced from the diagram associated to $f_{\mathbf{m}, \omega, \Sigma}$, for ϵ small enough. This bottleneck closeness occurs with large probability for a regular manifold \mathcal{X} for the c -PDTM built from a noisy sample from \mathcal{X} , according to [6]. No such result has been proved yet for the c -PLM. Anyway, intuitively, its sublevel sets are good approximations of the manifold \mathcal{X} , with the advantage that they are made of less ellipsoids, and that these ellipsoids are oriented accordingly to the manifold, i.e. with large eigenvalues on the tangent space and small eigenvalues on its orthogonal. This will be confirmed in the numerical illustrations section.

By construction, the persistence diagram (for connected components) associated to the filtration of the sublevel sets of $f_{\mathbf{m}, \omega, \Sigma}$ coincides with the persistence diagram associated to the anisotropic weighted Čech complex $\text{Cech}(f_{\mathbf{m}, \omega, \Sigma})$. Consequently, we can forget about the ellipsoids and focus on the simplicial complex filtration, which can be computed and stored efficiently, in a $c \times c$ matrix $\text{Mat} = (t_{i,j})_{i,j \in I}$. Such a matrix contains the times of appearance of vertices and of merging of connected components in $\text{Cech}(f_{\mathbf{m}, \omega, \Sigma})$. The clustering scheme of this paper exposed just below is based on such a merging matrix Mat .

3.2 An algorithm for persistence-based clustering

Consider $(\mathcal{G}^t)_{t \in \mathbb{R}}$ a filtration of sub-graphs of \mathcal{G} , a graph with c nodes. Based on this filtration, we define an algorithm, strongly inspired from the ToMATo algorithm [12]. The clustering scheme is guided by the persistence of the connected components in $(\mathcal{G}^t)_{t \in \mathbb{R}}$, and preserves components with large prominence. We assume that the nodes of \mathcal{G} are labeled such that the node labeled i gets born before the node labeled j , when $i \leq j$. The procedure is as follows. A connected component gets born when a node gets born, with the same label. A component changes of label at each time t for which it merges with a component with smaller label in \mathcal{G}^t , unless its prominence is larger than some parameter *Stop*. The prominence of a node or a component is defined as the lifetime of the component in the filtration (i.e. the elapsed time between the birth of the node and the time t such that a node with smaller index is present in its connected component in \mathcal{G}^t). The resulting clustering is given by the label of the nodes at time $t = +\infty$. It contains exactly labels of edges with a prominence larger than *Stop*. In this clustering scheme, we decide that nodes born after some time parameter *Threshold* are not relevant; they are removed. This procedure is implemented in Algorithm 1.

■ **Algorithm 1** Persistence-based Clustering Algorithm.

Data: Mat, Threshold, Stop
Result: Color, Birth, Death
Initialization ;
 $c \leftarrow \max\{i \mid \text{Mat}[i,i] \leq \text{Threshold}\}$; Mat \leftarrow Mat[1:c,1:c] ;
 Birth \leftarrow [Mat[i,i] for i in 1:c] ; Death \leftarrow [∞ for i in 1:c] ;
 indice \leftarrow 1 ; I \leftarrow 1 ; time \leftarrow Mat[I,I] ; Color \leftarrow [] ;
while time $<$ ∞ **do**
 if time = Mat[I,I] **then**
 Component I appears ;
 indice \leftarrow indice + 1 ; Mat[I,I] \leftarrow ∞ ; Color[I] \leftarrow I ;
 else
 (col_max, col_min) \leftarrow (max(Color[I],Color[J]) , min(Color[I],Color[J]));
 if time - Birth[col_max] \leq Stop **then**
 Components col_max and col_min merge ;
 Replace all entries col_max by col_min in Color ;
 Death[col_max] \leftarrow time ;
 else
 Component col_max will never die ;
 end
 Mat[i,j] \leftarrow ∞ for every i, j \leq min(indice,c) such that
 (Color[i],Color[j]) \in {(col_min, col_max), (col_max, col_min)} ;
 end
 I,J \leftarrow arg min_{i,j \leq min(indice,c)} Mat[i,j] ; time \leftarrow Mat[I,J]
end

This algorithm requires a merging matrix $\text{Mat} = (t_{i,j})_{i,j \in I}$, with $I = \llbracket 1, c \rrbracket$. We define its coefficients by $t_{i,i}$, the birth time of the node i in the filtration $(\mathcal{G}^t)_{t \in T}$; for $i > j$, $t_{i,j}$ the birth time of the edge $[i, j]$ and for $i < j$, $t_{i,j} = \infty$. The vector *Color* contains the resulting clustering, the vector *Birth*, the birth time of the components and *Death* their death time. Note that *Death*[1] is always $+\infty$. When $(\mathcal{G}_t)_{t \in T}$ is the filtration of the sublevel sets of some power function $f_{\mathbf{m}, \omega, \Sigma}$, the matrix Mat has coefficients given by $t_{i,i} = \omega_i$ and for $i > j \geq 1$, $t_{i,j}$ the intersecting time of the ellipsoids \mathcal{E}_i^t and \mathcal{E}_j^t , given by Proposition 1.

In practice, to label points in \mathbb{X} (generated around \mathcal{X}), we consider an approximation of $d_{\mathcal{X}}^2$ based on a family \mathbf{m} of c centers. Set \mathbf{m}' , the centers not removed and labeled by Algorithm 1, and $\boldsymbol{\omega}'$ and $\boldsymbol{\Sigma}'$ the corresponding parameters. Clustering points in \mathbb{X} is made accordingly to these labels and to the Voronoi decomposition of \mathbb{R}^d , based on \mathbf{m}' , $\boldsymbol{\omega}'$ and $\boldsymbol{\Sigma}'$: $x \in \mathbb{X}$ has the same label as m'_i if $\|x - m'_i\|_{\boldsymbol{\Sigma}'_i}^2 + \omega'_i \leq \|x - m'_j\|_{\boldsymbol{\Sigma}'_j}^2 + \omega'_j$ for every j . Since $f_{\mathbf{m}', \boldsymbol{\omega}', \boldsymbol{\Sigma}'}$ approximates $d_{\mathcal{X}}^2$, in order to deal with outliers, we remove (i.e. assign the label 0) the points $x \in \mathbb{X}$ for which $f_{\mathbf{m}', \boldsymbol{\omega}', \boldsymbol{\Sigma}'}(x)$ is the largest. Note that a power function is homogeneous to the square of a distance function. Therefore, for positive weights $\boldsymbol{\omega}$, it could be more appropriate to consider the filtration of sublevel sets of $\sqrt{f_{\mathbf{m}, \boldsymbol{\omega}, \boldsymbol{\Sigma}}}$ instead of $f_{\mathbf{m}, \boldsymbol{\omega}, \boldsymbol{\Sigma}}$.

The best complexity of Algorithm 1 ($O(c^3)$ comparisons) is obtained when $Stop = \infty$, with $2c$ iterations of the algorithm. The worst complexity ($O(c^4)$) is obtained when $Stop = 0$, with $O(c^2)$ iterations. This is fast when c is much smaller than the sample size (e.g. for c -PLM and c -PDTM), and does not depend on the dimension. In the experiments of Section 4, Algorithm 1 runs much faster than the computation of the c -PLM and the c -PDTM.

In practice, just as Chazal et al. [12], we recommend to run Algorithm 1 several times. A first time with $Threshold = Stop = \infty$ to calibrate the parameter $Threshold$, in order to remove bad nodes (i.e. nodes with late birth and short lifetime). A second time with this parameter $Threshold$ and $Stop = \infty$, to measure the prominence of the components and select the number of clusters (via the parameter $Stop$), as the number of components with prominence much larger than others. More details on the calibration of these two parameters, from the persistence diagrams $(Birth[i], Death[i])_{i \in I}$, are given in Section 4.1. The final clustering is obtained from $Color$, after running Algorithm 1 with these two parameters.

Giving a sense to an optimal minimal prominence $Stop$ is possible for distance functions. For instance, for the sublevel-sets filtration of $d_{\mathcal{X}}$, $Stop$ can be chosen as half of the minimal distance between two distinct components of \mathcal{X} . Consequently, for any $\epsilon - \|\cdot\|_{\infty}$ -close approximation of $d_{\mathcal{X}}$, taking $Stop - \epsilon$ leads to a perfect clustering, provided that $2\epsilon < Stop$.

The parameter $Threshold$ is primordial, especially for the c -PLM function. Indeed, the algorithm for the c -PLM is based on $\tilde{\theta}_{c,k}$, a local minimizer of the criterion $R_{c,k}$. Consequently, some ellipsoids \mathcal{E}_i are far from the support, or in a wrong direction. Thus, their weight ω_i (and thus $Birth[i]$) is large with respect to other well-placed ellipsoids, due to a large variance term $v_{y_i, \Sigma_i, k}$. Such bad ellipsoids are removed for a suitable parameter $Threshold$.

3.3 Connection to other persistence-based clustering methods

In the sequel, we display different graph filtrations, to be used for persistence-based clustering, with Algorithm 1. For each of these filtrations, we give a summarize of the corresponding matrices Mat , in Table 1, with the convention that $t_{i,i} \leq t_{j,j}$ when $i \leq j$.

ToMATo Algorithm [12] rests on a graph filtration based on a graph \mathcal{G} and a function f defined on the nodes of \mathcal{G} . Morally, \mathcal{G}^t is the sub-graph of \mathcal{G} that contains the nodes i such that $f(i) \leq t$, and the edges $[i, j]$ if and only if i and j are in \mathcal{G}^t . Chazal et al. mostly studied this method for \mathcal{G} , a Rips graph of a set $\mathbb{X} \subset \mathbb{R}^d$, and for $f(i)$, the DTM to \mathbb{X} at X_i .

The DTM-filtration [1] corresponds to the 1-skeleton of the nerve of the union of balls $(\bigcup_{x \in \mathbb{X}} \bar{B}(x, r_t(x)))_{t > 0}$ with $r_t(x) = -\infty$ for $t < d_{\mathbb{X}, k}(x)$ and $r_t(x) = (t^p - d_{\mathbb{X}, k}^p(x))^{\frac{1}{p}}$ for $t \geq d_{\mathbb{X}, k}(x)$, for some $p \geq 1$ and with the convention that $\bar{B}(x, -\infty)$ is empty. In Table 1, we give the coefficients for $p = 1$. The DTM-filtration with $p = 2$ was actually introduced in [7], leading to what we call Power filtration, which coincides with the sublevel-sets filtration of the square of a power distance. We also consider additional power-functions-based filtrations, from the k -witnessed distance [18], the c -PDTM [6] and the c -PLM [4].

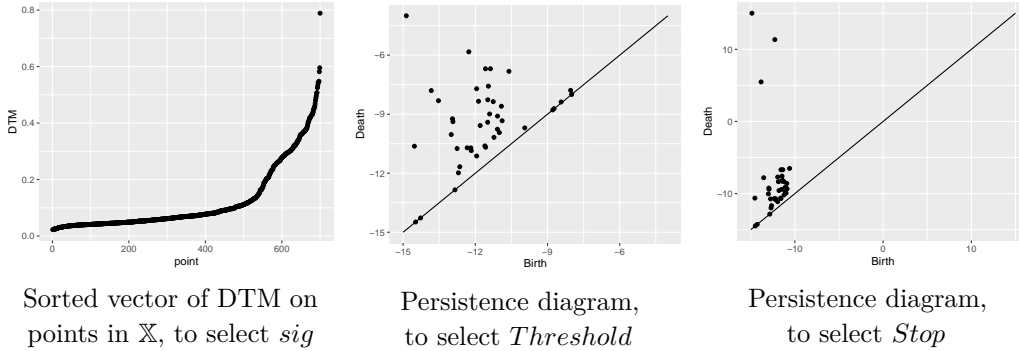
■ **Table 1** Coefficients of Mat for the different methods, with the notation $f = d_{\mathbb{X},k}$ for the DTM to \mathbb{X} with number of nearest neighbors parameter k .

Method	$t_{i,i}$	$t_{i,j}$ for $i < j$
ToMATo	$f(i)$	$\max(f(i), f(j))(\mathbb{1}_{[i,j] \in \mathcal{G}})^{-1}$
DTM-filtration	$f(i)$	$\left(\frac{\ X_i - X_j\ + f(i) + f(j)}{2} \right) \mathbb{1}_{\ X_i - X_j\ > f(i) - f(j) } + f(i) \mathbb{1}_{f(i) - f(j) \geq \ X_i - X_j\ }$
$f_{\mathbf{m},\omega}$	ω_i	$\frac{(\omega_j - \omega_i)^2 + 2(\omega_j + \omega_i)\ m_j - m_i\ ^2 + \ m_j - m_i\ ^4}{4\ m_j - m_i\ ^2}$
$\sqrt{f_{\mathbf{m},\omega}}$	$\sqrt{\omega_i}$	$\sqrt{\frac{(\omega_j - \omega_i)^2 + 2(\omega_j + \omega_i)\ m_j - m_i\ ^2 + \ m_j - m_i\ ^4}{4\ m_j - m_i\ ^2}}$
$f_{\mathbf{m},\omega,\Sigma}$	ω_i	Given by Proposition 1
Power filtration		$\sqrt{f_{\mathbf{m},\omega}}$ with $\mathbf{m} = \mathbb{X}$ and $\omega = (f^2(x))_{x \in \mathbb{X}}$
Witnessed		$\sqrt{f_{\mathbf{m},\omega}}$ with $(\mathbf{m}, \omega) = (m_{x,k}, v_{x,k})_{x \in \mathbb{X}}$
c -PDTM		$f_{\mathbf{m},\omega}$ with $(\mathbf{m}, \omega) = (m_{y,k}, v_{y,k})_{y \in \mathcal{Y}_{c,k}}$
c -PLM		$f_{\mathbf{m},\omega,\Sigma}$ with $(\mathbf{m}, \omega, \Sigma) = (m_{y,\Sigma,k}, v_{y,\Sigma,k} + \log(\det(\Sigma)), \Sigma)_{(y,\Sigma) \in \theta_{c,k}}$

4 Numerical illustrations

4.1 A complete illustration of the method

Consider the target \mathcal{X} , a set of three curves in \mathbb{R}^2 . We generate $\mathbb{X} = (X_i)_{i \in [1, N_s + N_o]}$, a set of $N_s = 500$ signal points $(X_i = Y_i + Z_i)_{i \in [1, N_s]}$, with Y_i uniform on \mathcal{X} and Z_i Gaussian with standard deviation $\sigma = 0.02$; corrupted by $N_o = 200$ outliers, uniform on $[-1.5, 2.5]^2$. We compare the clustering scheme based on Algorithm 1 with the sublevel sets of the c -PLM, to the target labels in Figure 2 (left). Parameters are set to $c = 50$ centers, $k = 10$ nearest neighbors, $sig = 520$ points to consider as signal, and $it = 100$ iterations and $n_ini = 10$ initializations to compute a suitable local optimum $\tilde{\theta}_{c,k,sig}$ of the c -PLM-criterion $R_{c,k,sig}$. Since the DTM $d_{\mathbb{X},k}$ is large for outliers, we select sig from the curve $([d_{\mathbb{X},k}(X_i), i \in [1, N_s + N_o]]$ in non-decreasing order), as the point of slope break; see Figure 1 (left). The DTM can be replaced by any not-trimmed approximation of the c -PLM.



■ **Figure 1** Parameters selection heuristics.

We run Algorithm 1 a first time with the parameters $Threshold = \infty$ and $Stop = \infty$, and display the persistence diagram $(Birth[i], Death[i])_{i \in [1, c]}$, in Figure 1 (middle). In order to have 3 clusters, we select $Stop = 5.62$, the height of a line parallel to the diagonal, separating 3 points from the others. We run Algorithm 1 a second time with this new parameter, which results in the clustering \mathcal{C}_1 of Figure 2 (middle). A sublevel set of the function $f_{\tilde{\theta}_{c,k}}$ is represented by the union of ellipses. Note that some ellipses have a bad position. This results

in a bad clustering. We use the parameter *Threshold* to remove them. In Figure 1 (middle), 6 points are on the right side, separated from the other points with a vertical line (of abscissa -10.27). Then, we run Algorithm 1 with $Threshold = -10.27$ and $Stop = \infty$. According to the persistence diagram in Figure 1 (right), since 3 points are well-separated from the other ones with a large band parallel to the diagonal (containing a line parallel to the diagonal, with height 12), we recover the number of clusters, 3, and set $Stop = 12$. The clustering \mathcal{C}_2 obtained with $Threshold = -10.27$ and $Stop = 12$ is represented in Figure 2 (right). The bad ellipses have been removed. Denote by $\tilde{\theta}'_{c,k,sig}$, the subfamily of $\tilde{\theta}_{c,k,sig}$ made of centers not removed by the procedure. The color of any point x in Figure 2 (right) is given by the label in *Color* (label returned by the Algorithm 1) of its associated center (y, Σ) in $\tilde{\theta}'_{c,k,sig}$. This is the center (y, Σ) such that $f_{\tilde{\theta}'_{c,k,sig}}(x) = \|x - m_{y,\Sigma,k}\|_{\Sigma^{-1}}^2 + v_{y,\Sigma,k} + \log(\det(\Sigma))$. The labels of the $|\mathbb{X}| - sig$ points with largest $f_{\tilde{\theta}'_{c,k,sig}}$ -value are set to 0.

Note that for large datasets, computing $\tilde{\theta}'_{c,k,sig}$ may take some time. We can compute it from a sub-sample of \mathbb{X} , run Algorithm 1, and label points in \mathbb{X} accordingly.

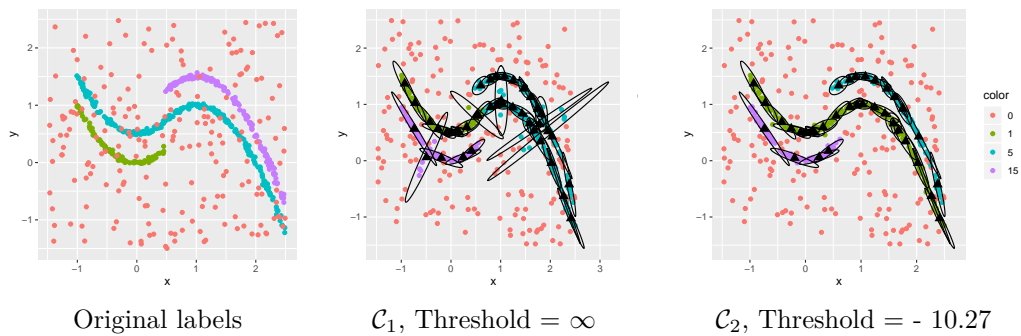


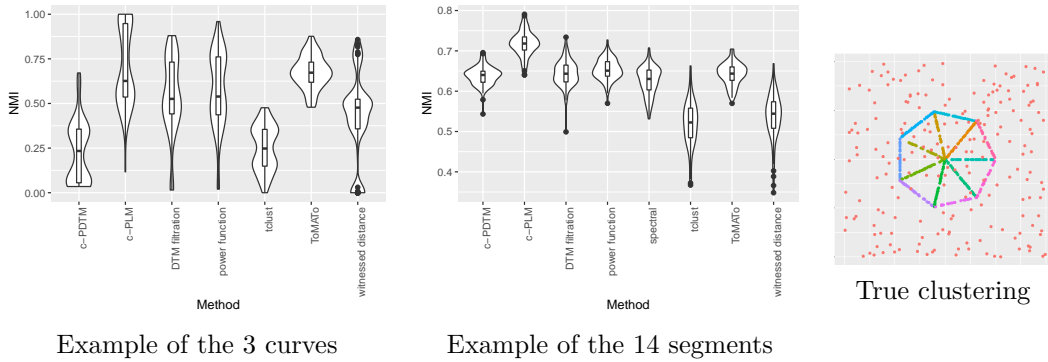
Figure 2 Two resulting clusterings, with ellipses.

We compare the performance of the two clusterings \mathcal{C}_1 and \mathcal{C}_2 . In terms of outliers detection, this can be assessed via the proportion of signal points labeled as outliers (0.034 for \mathcal{C}_1 , 0.016 for \mathcal{C}_2) and as the proportion of outliers labeled as signal points (0.185 for \mathcal{C}_1 , 0.14 for \mathcal{C}_2). As expected from Figure 2, removing bad ellipses reduces these proportions and thus improves the outliers detection performance. In terms of clusters recovering, the normalized mutual information (NMI) is classically used. It equals 1 for a perfect clustering and 0 for a terrible clustering. When considering outliers as a cluster with label 0, we got $NMI = 0.586$ for \mathcal{C}_1 and $NMI = 0.841$ for \mathcal{C}_2 . The NMI computed on the signal points labeled as signal points is $NMI = 0.634$ for \mathcal{C}_1 and $NMI = 1$ for \mathcal{C}_2 , a perfect clustering.

4.2 Comparison of the different methods on synthetic datasets

We compare different clustering methods on two synthetic datasets : the previous dataset with 3 curves, and datapoints from a polygonal curve of 14 segments, as in [8]. We set parameters to $N_s = 500$, $N_o = 200$, $\sigma = 0.02$, $c = 50$, $k = 10$, $it = 100$, $n_{ini} = 10$ and *Threshold* chosen such that 10 means are removed from the c -PLM-centers $\tilde{\theta}_{c,k,sig}$. For the ToMATo algorithm we set $r = 0.12$, the radius of the Rips graph. We used the function *dbscan* from the R packages *dbscan* [19], with parameters $eps = 0.15$ and $minPts = 10$; *tclust* and *specc* from the *tclust* [17] and *kernlab* [21] R packages.

For the three curves, the parameter r for ToMATo is chosen such that the graph is not connected, the clusterings are acceptable but have more than 3 clusters. The c -PLM often performs perfectly, and sometimes performs poorly, since the number of bad ellipses removed



■ **Figure 3** Violin plots representing the NMI computed on signal points, detected as signal points.

is fixed to 10 and not calibrated according to the heuristics, and there is some instability. We observe the same clustering problem as in Figure 2 (middle) for the other methods since the lines are close, compared to the distance between sample points from the same line. For the polygonal line of 14 segments, all methods except the *c*-PLM and *tclust* put centers of clusters on massive parts of \mathbb{X} (the center and the intersections of 3 segments). For the *c*-PLM and *tclust*, most clusters coincide with segments. Nonetheless, there is some instability (much less pronounced for the *c*-PLM), since the algorithms are based on local minimizers.

4.3 Applications to real datasets

4.3.1 Recovering fleas species, based on 6 measurements

We picked the dataset flea from the R-package *tour* [29], initially from [23]. This dataset contains records of 6 measurements for 74 male insects from the Palaearctic, from three different species: *Heptapotamica*, *Concinna*, *Heikertingeri*. The variables correspond to measurements on the tarsus, the aedeagus and the head. We normalized data so that the mean and variance of each of the 6 variables are respectively 0 and 1. In Table 2, we computed the NMI between the true species and the clustering returned by different methods. We ran each algorithm 10 times with at most 100 iterations. For every *k*-nearest-neighbours-based algorithm, we set $k = 10$. For ToMATo, we set $r = 1.9$ so that the graph is connected; for the *c*-PLM and the *c*-PDTM, $c = 50$ and for *dbscan*, $eps = 1.5$ and $minPts = 10$. The 3-PDTM and 3-PLM methods consist in clustering data according to the weighted Voronoi cells given by the optimal centers and covariance matrices.

■ **Table 2** NMI between clustering of fleas and their true species.

Without Algorithm 1	<i>k</i> -means 0.825	<i>tclust</i> 0.848	DBSCAN 0.647	Spectral 1	3-PLM 1	3-PDTM 1
With Algorithm 1	ToMATo 0.628	Witnessed 0.906	power 1	DTM-filt. 1	<i>c</i> -PLM hier. 1	<i>c</i> -PDTM hier. 1

The methods based on the decomposition of \mathbb{R}^6 into 3 (weighted and/or curved) Voronoi cells are efficient: at most 3 bad labels for *k*-means and *tclust* and all labels correct for their “robust” versions, the 3-PDTM and the 3-PLM. The perfect performance of these two last functions is due to the weights that force the centers of cells to lie in massive areas for \mathbb{X} . The bad performance of ToMATo is due to the difficulty to select the parameter r for the Rips

graph, the small number of points, and the fact that the inverse of the DTM should be used instead of the DTM, as recommended by the authors. Nonetheless, we made the choice to use the DTM since the other methods (witnessed distance, power function, DTM-filtration, c -PLM and c -PDTM) are based on filtrations from approximations of the DTM, and almost all of these methods perform perfectly. The method `dbscan` performs poorly since it labels 14 points as outliers. Nonetheless, the points considered as signal are well clustered.

4.3.2 Clustering a earthquake dataset

We consider a set of 12790 points representing the longitude and latitude of earthquakes of magnitude non smaller than 5.0, between the 01/01/1970 and the 01/01/2010. This dataset was picked from the website <http://earthquake.usgs.gov/earthquakes/eqarchives/epic/>.

We used Algorithm 1 with an approximation of the c -PLM based on a sub-sample of 2000 points from the dataset, with parameters $c = 200$, $k = 10$ and for $it = 50$ iterations. We restricted matrices Σ to have eigenvalues smaller than 50 by thresholding them. The persistence diagram in Figure 4 suggests that the dataset has 4 or 10 clusters. Moreover, the curve of the sorted values of the c -PLM approximation on the pointset in Figure 4 suggests to keep $sig = 12250$ points as signal points. See Figure 5 for the corresponding clustering.

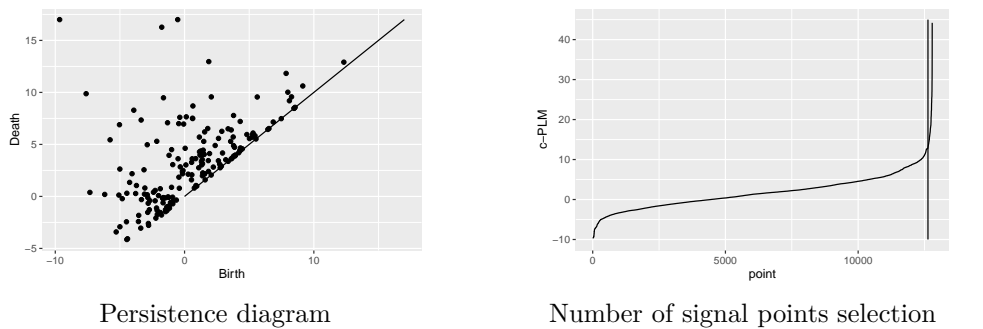


Figure 4 Parameters selection heuristics.

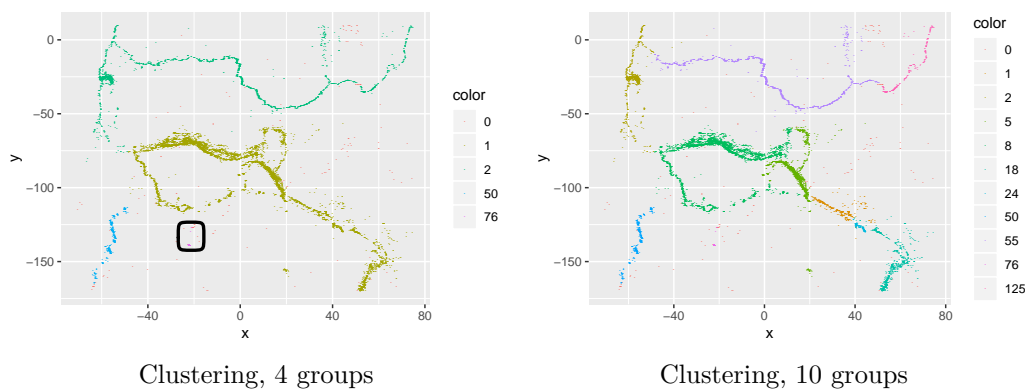


Figure 5 Earthquake clustering with Algorithm 1, for the c -PLM function.

References

- 1 Hirokazu Anai, Frédéric Chazal, Marc Glisse, Yuichi Ike, Hiroya Inakoshi, Raphaël Tinarrage, and Yuhei Umeda. DTM-based filtrations. In *35th International Symposium on Computational Geometry*, volume 129 of *LIPICs. Leibniz Int. Proc. Inform.*, pages Art. No. 58, 15. Schloss Dagstuhl. Leibniz-Zent. Inform., Wadern, 2019.
- 2 Arindam Banerjee, Srujana Merugu, Inderjit S. Dhillon, and Joydeep Ghosh. Clustering with bregman divergences. *J. Mach. Learn. Res.*, 6:1705–1749, December 2005. URL: <http://dl.acm.org/citation.cfm?id=1046920.1194902>.
- 3 Gregory Bell, Austin Lawson, Joshua Martin, James Rudzinski, and Clifford Smyth. Weighted persistent homology. *Involve*, 12(5):823–837, 2019. doi:10.2140/involve.2019.12.823.
- 4 Claire Bréchet. Robust shape inference from a sparse approximation of the gaussian trimmed loglikelihood. Unpublished, 2018.
- 5 Claire Bréchet, Aurélie Fischer, and Clément Levrard. Robust bregman clustering. In revision, 2018.
- 6 Claire Bréchet and Clément Levrard. A k-points-based distance for robust geometric inference. To appear in Bernoulli, 2017.
- 7 Mickaël Buchet, Frédéric Chazal, Steve Y. Oudot, and Donald R. Sheehy. Efficient and robust persistent homology for measures. *Comput. Geom.*, 58:70–96, 2016. doi:10.1016/j.comgeo.2016.07.001.
- 8 Mickaël Buchet, Tamal K. Dey, Jiayuan Wang, and Yusu Wang. Declutter and resample: towards parameter free denoising. *J. Comput. Geom.*, 9(2):21–46, 2018.
- 9 Frédéric Chazal, David Cohen-Steiner, Marc Glisse, Leonidas J. Guibas, and Steve Y. Oudot. Proximity of persistence modules and their diagrams. In *Proceedings of the Twenty-fifth Annual Symposium on Computational Geometry, SCG '09*, pages 237–246, New York, NY, USA, 2009. ACM. doi:10.1145/1542362.1542407.
- 10 Frédéric Chazal, David Cohen-Steiner, and Quentin Mérigot. Geometric Inference for Measures based on Distance Functions. *Foundations of Computational Mathematics*, 11(6):733–751, 2011. doi:10.1007/s10208-011-9098-0.
- 11 Frédéric Chazal, Vin de Silva, Marc Glisse, and Steve Oudot. *The structure and stability of persistence modules*. SpringerBriefs in Mathematics. Springer, [Cham], 2016. doi:10.1007/978-3-319-42545-0.
- 12 Frédéric Chazal, Leonidas J. Guibas, Steve Y. Oudot, and Primoz Skraba. Persistence-based clustering in Riemannian manifolds. *J. ACM*, 60(6):Art. 41, 38, 2013. doi:10.1145/2535927.
- 13 David Cohen-Steiner, Herbert Edelsbrunner, and John Harer. Stability of persistence diagrams. *Discrete Comput. Geom.*, 37(1):103–120, 2007. doi:10.1007/s00454-006-1276-5.
- 14 J. A. Cuesta-Albertos, A. Gordaliza, and C. Matrán. Trimmed k -means: an attempt to robustify quantizers. *Ann. Statist.*, 25(2):553–576, 1997. doi:10.1214/aos/1031833664.
- 15 Vin de Silva and Robert Ghrist. Coverage in sensor networks via persistent homology. *Algebr. Geom. Topol.*, 7:339–358, 2007. doi:10.2140/agt.2007.7.339.
- 16 Herbert Edelsbrunner, David Letscher, and Afra Zomorodian. Topological persistence and simplification. *Discrete Comput. Geom.*, 28(4):511–533, 2002. Discrete and computational geometry and graph drawing (Columbia, SC, 2001). doi:10.1007/s00454-002-2885-2.
- 17 Heinrich Fritz, Luis A. Garcia-Escudero, and Agustin Mayo-Iscar. tclust: An R package for a trimming approach to cluster analysis. *Journal of Statistical Software*, 47(12):1–26, 2012. URL: <http://www.jstatsoft.org/v47/i12/>.
- 18 Leonidas Guibas, Dmitriy Morozov, and Quentin Mérigot. Witnessed k -distance. *Discrete Comput. Geom.*, 49(1):22–45, 2013. doi:10.1007/s00454-012-9465-x.
- 19 Michael Hahsler, Matthew Piekenbrock, and Derek Doran. dbscan: Fast density-based clustering with R. *Journal of Statistical Software*, 91(1):1–30, 2019. doi:10.18637/jss.v091.i01.
- 20 Allen Hatcher. *Algebraic topology*. Cambridge University Press, Cambridge, 2002.

- 21 Alexandros Karatzoglou, Alex Smola, Kurt Hornik, and Achim Zeileis. kernlab – an S4 package for kernel methods in R. *Journal of Statistical Software*, 11(9):1–20, 2004. URL: <http://www.jstatsoft.org/v11/i09/>.
- 22 Stuart P. Lloyd. Least squares quantization in PCM. *IEEE Trans. Inform. Theory*, 28(2):129–137, 1982. doi:10.1109/TIT.1982.1056489.
- 23 Alexander A. Lubischew. On the use of discriminant functions in taxonomy. *Biometrics*, pages 455–477, 1962.
- 24 J. MacQueen. Some methods for classification and analysis of multivariate observations. In *Proceedings of the Fifth Berkeley Symposium on Mathematical Statistics and Probability, Volume 1: Statistics*, pages 281–297, Berkeley, Calif., 1967. University of California Press. URL: <https://projecteuclid.org/euclid.bsmsp/1200512992>.
- 25 Stephen B Pope. Algorithms for ellipsoids. Technical Report FDA-08-01, Sibley School of Mechanical & Aerospace Engineering, Cornell University Ithaca, New York 14853, 2008.
- 26 P. J. Rousseeuw and A. M. Leroy. *Robust Regression and Outlier Detection*. John Wiley & Sons, New York, 1987.
- 27 Ulrike von Luxburg. A tutorial on spectral clustering. *Stat. Comput.*, 17(4):395–416, 2007. doi:10.1007/s11222-007-9033-z.
- 28 Wenping Wang, Jiaye Wang, and Myung-Soo Kim. An algebraic condition for the separation of two ellipsoids. *Comput. Aided Geom. Design*, 18(6):531–539, 2001. doi:10.1016/S0167-8396(01)00049-8.
- 29 Hadley Wickham, Dianne Cook, Heike Hofmann, and Andreas Buja. tourr: An R package for exploring multivariate data with projections. *Journal of Statistical Software*, 40(2):1–18, 2011. URL: <http://www.jstatsoft.org/v40/i02/>.