

Robust Anisotropic Power-Functions-based Filtrations for Clustering.

Claire BréchetEAU 

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

1 Abstract

2 We consider robust power-distance functions that approximate the distance function to a compact
3 set, from a noisy sample. We pay particular interest to robust power-distance functions that are
4 anisotropic, in the sense that their sublevel sets are unions of ellipsoids, and not necessarily unions
5 of balls. Using persistence homology on such power-distance functions provides robust clustering
6 schemes. We investigate such clustering schemes and compare the different procedures on synthetic
7 and real datasets. In particular, we enhance the good performance of the anisotropic method for
8 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

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

Supplement 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.

Lines 476

9 1 Introduction

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

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



© C. BréchetEAU;

licensed under Creative Commons License CC-BY

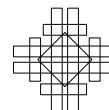
36th International Symposium on Computational Geometry (SoCG 2020).

Editors: Sergio Cabello and Danny Z. Chen; Article No. 23; pp. 23:1–23:15

Leibniz International Proceedings in Informatics



LIPIC Schloss Dagstuhl – Leibniz-Zentrum für Informatik, Dagstuhl Publishing, Germany



30 sets of the distance function to \mathbb{X} . For a fixed radius r , the Čech complex is a simplicial
 31 complex defined as the collection of simplices (vertex, edge, triangle, tetrahedron) for which
 32 the r -balls centered at the vertices have a non-empty common intersection. We call 1-skeleton
 33 its subcomplex (a graph) that contains only vertices and edges. The non-decreasing family
 34 of such graphs indexed by $r \in \mathbb{R}$ is called a filtration. Single-linkage is a persistence-based
 35 method since is based on the persistence, prominence or equivalently lifetime of the connected
 36 components into this graph filtration, however, it is not robust to outliers. The algorithm
 37 ToMATo in [12] is robust and persistence-based. Indeed, it is based on a graph filtration
 38 built from a neighborhood graph and a (robust) distance-like function whose values guide
 39 the appearance of vertices and edges in the graph filtration. An example of robust distance
 40 function that Chazal et al. consider in [12] is given by the distance-to-measure (DTM) [10].
 41 Note that the graph is a priori not intrinsic to the distance function, which may cause bad
 42 clustering. For instance, an edge that links two vertices with small distance-function value
 43 but intersects an area with large distance function value, may link two clusters that should
 44 not be. This problem was the cause of failure of the single-linkage method for data corrupted
 45 by outliers. Alternative filtrations that do not suffer from this problem are the DTM-filtration
 46 [1], or the power filtrations [7], based on the 1-skeleton of the Čech filtration associated to
 47 the sublevel sets of a power distance function: a function of type $x \mapsto \min_{i \in I} \|x - m_i\|^2 + \omega_i$
 48 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
 49 functions have been introduced and studied in the literature: the k -witnessed distance [18],
 50 the power distance [7], the c -PDTM [6] whose sublevel sets are unions of c balls, and the
 51 c -PLM [4] whose sublevel sets are unions of c ellipsoids, with c possibly much smaller than
 52 the sample size. The last two functions are robust to outliers since their construction is
 53 based on the principle of trimmed least squares [26].

54 Contributions

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

61 Practical interests

62 A clustering algorithm based on the persistence filtration of the sublevel sets of a power
 63 function is pertinent since unlike ToMATo, the graph is intrinsic to the distance function.
 64 So, no additional parameters are required for the algorithm. The main advantage of using
 65 an anisotropic power function is that its sublevel sets are ellipsoids. Much less ellipsoids are
 66 required than balls to Hausdorff-approximate a compact manifold with intrinsic dimension
 67 smaller than the ambient dimension. The clustering scheme can also be applied to decompose
 68 a set of points generated on a polygonal line into segments. Once the ellipsoids computed,
 69 the persistence algorithm runs fast. Its complexity in terms of number of comparisons is at
 70 worst $O(c^4)$, with c , the number of ellipsoids, which is much smaller than the sample size.
 71 Most importantly, the robustness of the persistence algorithm relies on the robustness of
 72 the distance function. The c -PLM [4] is robust to outliers. The guaranty for the clustering
 73 method follows from the $\|\cdot\|_\infty$ -distance closeness between the power distance function and the
 74 distance function to the underlying manifold \mathcal{X} , relatively to the minimal distance between

75 the connected components of \mathcal{X} . Note that such a proximity condition is sufficient but not
 76 necessary, as illustrated by the different numerical examples, with the c -PLM.

77 Organisation of the paper

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

87 **2 Power-functions-based filtrations for robust clustering**

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

97 2.1 Generalities on filtrations

98 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
 99 for the inclusion (i.e. $\forall t \leq t', V^t \subset V^{t'}$). A typical example is the filtration of the sub-level
 100 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} ,
 101 a filtration of simplicial complexes of S is a non-decreasing family $(S^t)_{t \in T}$ of subcomplexes
 102 of S , meaning that for every $t \leq t'$, any simplex of S^t is also a simplex of $S^{t'}$.

103 The interleaving pseudo-distance between two filtrations $(V^t)_{t \in T}$ and $(W^t)_{t \in T}$ is defined
 104 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:
 105 $\forall t \in T, V^t \subset W^{t+\epsilon}$ and $W^t \subset V^{t+\epsilon}$. This definition extends to simplicial complexes. Note
 106 that the sub-level-sets filtrations of two functions f and g satisfying $\|f - g\|_\infty \leq \epsilon$ are
 107 ϵ -interleaved. We will see in Section 3 that the notion of interleaving is primordial, since it
 108 measures the difference of topology between two filtrations. In particular, the stability of our
 109 sub-level-sets-based clustering scheme will be guaranteed from the closeness of the functions.

110 2.2 Power-functions-based filtrations

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

115 **2.2.1 The isotropic case**

116 An isotropic power function is a function $f_{\mathbf{m},\omega} : \mathbb{R}^d \rightarrow \mathbb{R}$ defined from an index set $I = \llbracket 1, c \rrbracket$,
 117 a family of centers $\mathbf{m} = (m_i)_{i \in I}$ in \mathbb{R}^d and a family of weights $\omega = (\omega_i)_{i \in I}$ in \mathbb{R} by
 118 $f_{\mathbf{m},\omega} : x \mapsto \min_{i \in I} \|x - m_i\|^2 + \omega_i$. A simple example of power function is the squared
 119 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
 120 sublevel sets of $f_{\mathbf{m},\omega}$, $V_{\mathbf{m},\omega}^t = f_{\mathbf{m},\omega}^{-1}((-\infty, t])$, are unions of at most c balls $\mathcal{B}_i^t = \overline{B}(m_i, \sqrt{t - \omega_i})$
 121 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
 122 \mathcal{B}_i^t and \mathcal{B}_j^t intersect if and only if $t \geq t_{i,j}$ with $t_{i,j} = \frac{(\omega_j - \omega_i)^2 + 2(\omega_j + \omega_i)\|m_j - m_i\|^2 + \|m_j - m_i\|^4}{4\|m_j - m_i\|^2}$.
 123 The connectivity of $V_{\mathbf{m},\omega}^t$ can be encoded in a graph $\mathcal{G}_{\mathbf{m},\omega}^t$, whose vertices are indices $i \in I$
 124 such that $\omega_i \leq t$ and whose edges are pairs of vertices $[i, j]$ such that $t_{i,j} \leq t$. Indeed, $\mathcal{G}_{\mathbf{m},\omega}^t$
 125 and $V_{\mathbf{m},\omega}^t$ have the same number of connected components, and m_i and m_j are in the same
 126 connected component in $V_{\mathbf{m},\omega}^t$ if and only if i and j are also in the same component in $\mathcal{G}_{\mathbf{m},\omega}^t$.

127 More generally, the topological information of $V_{\mathbf{m},\omega}^t$ (number of connected components,
 128 loops, voids etc.) can be encoded in the weighted Čech complex $\text{Cech}_{\mathbf{m},\omega}(t)$, defined as
 129 the nerve of the union of balls $(\mathcal{B}_i^t)_{i \in I}$: $\text{Cech}_{\mathbf{m},\omega}(t) = \{\sigma \subset I \mid \bigcap_{i \in \sigma} \mathcal{B}_i^t \neq \emptyset\}$, [1, 7, 3].
 130 According to the Nerve Lemma [20, Corollary 4G.3], any sublevel set $V_{\mathbf{m},\omega}^t$ is homotopic
 131 to $\text{Cech}_{\mathbf{m},\omega}(t)$ and thus contains the same topological information. For computational
 132 reasons, the weighted Vietoris-Rips filtration is frequently considered as a provably good
 133 surrogate for the weighted Čech filtration $(\text{Cech}_{\mathbf{m},\omega}(t))_{t \in T}$. The weighted Vietoris-Rips
 134 complex $\text{VR}_{\mathbf{m},\omega}(t)$ is the flag complex of $\mathcal{G}_{\mathbf{m},\omega}^t$ ($\mathcal{G}_{\mathbf{m},\omega}^t$ is the 1-skeleton of the weighted Čech
 135 complex): $\text{VR}_{\mathbf{m},\omega}(t) = \{\sigma \subset I \mid \forall i, j \in \sigma, \mathcal{B}_i^t \cap \mathcal{B}_j^t \neq \emptyset\}$. Indeed, as a direct consequence of
 136 [3, Theorem 3.2] which is a generalization of the non-weighted case in [15, Theorem 2.5.], if
 137 the weights in ω are non-negative, then these two filtrations are interleaved:

$$138 \quad \forall 0 < t' \leq \frac{d+1}{2d}t, \text{VR}_{\mathbf{m},\omega}(t') \subset \text{Cech}_{\mathbf{m},\omega}(t) \subset \text{VR}_{\mathbf{m},\omega}(t). \quad (1)$$

139 These notions can all be extended to anisotropic power functions.

 140 **2.2.2 The anisotropic case**

141 Consider $I = \llbracket 1, c \rrbracket$, centers $\mathbf{m} = (m_i)_{i \in I}$ in \mathbb{R}^d , weights $\omega = (\omega_i)_{i \in I}$ in \mathbb{R} and matrices
 142 $\Sigma = (\Sigma_i)_{i \in I}$ in \mathcal{M}_d , the set of definite positive symmetric matrices. An anisotropic power
 143 function is a function $f_{\mathbf{m},\omega,\Sigma} : \mathbb{R}^d \rightarrow \mathbb{R}$ defined from I , \mathbf{m} , ω and Σ by $f_{\mathbf{m},\omega,\Sigma} : x \mapsto$
 144 $\min_{i \in I} \|x - m_i\|_{\Sigma_i}^2 + \omega_i$. For any matrix $\Sigma \in \mathcal{M}_d$ and $x \in \mathbb{R}^d$, $\|x\|_{\Sigma^{-1}} = \sqrt{x^T \Sigma^{-1} x}$ denotes
 145 the Σ -Mahalanobis norm of x . The sublevel sets of $f_{\mathbf{m},\omega,\Sigma}$, $V_{\mathbf{m},\omega,\Sigma}^t = f_{\mathbf{m},\omega,\Sigma}^{-1}((-\infty, t])$, are
 146 unions of at most c ellipsoids $\mathcal{E}_i^t = \overline{B}_{\Sigma_i}(m_i, \sqrt{t - \omega_i}) = \{x \in \mathbb{R}^d \mid \|x - m_i\|_{\Sigma_i}^2 \leq t - \omega_i\}$.
 147 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
 148 relative question of the emptiness of the intersection of two ellipsoids is tackled in [28, 25].

149 **► Proposition 1.** *Consider two ellipsoids $\mathcal{E}_i^t = \overline{B}_{\Sigma_i}(m_i, \sqrt{t - \omega_i})$ and $\mathcal{E}_j^t = \overline{B}_{\Sigma_j}(m_j, \sqrt{t - \omega_j})$
 150 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
 151 positive diagonal matrices D_i and D_j and two orthogonal matrices P_i and P_j from the spectral
 152 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
 153 $\tilde{D} = \text{diag}(\lambda_1, \lambda_2, \dots, \lambda_d)$, 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
 154 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 $t_{i,j} = \omega_j + \sum_{k=1}^d \left(\frac{\lambda_k \tilde{m}_k}{\lambda + \lambda_k} \right)^2 \lambda_k$
 155 when $\|\tilde{m}\| > \sqrt{\omega_j - \omega_i}$. The positive number λ is the unique solution of the following equation:*

156

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

158 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
 159 point, and the corresponding gradients are collinear. In the context of isotropy (i.e. for
 160 $\Sigma_i = \Sigma_j = I_d$, the identity matrix of \mathbb{R}^d) $\tilde{m} = m_j - m_i$, and when $\|m_j - m_i\| > \sqrt{\omega_j - \omega_i}$,
 161 (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
 162 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
 163 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,
 164 since unions of ellipsoids are contractible. Although this paper is mostly based on the study of
 165 connected components for clustering, anisotropic weighted Čech and Vietoris-Rips filtrations
 166 are primordial to have a tractable estimation of the topology of compact sets from suitable
 167 approximations as finite unions of ellipsoids. In fact, as their isotropic counterparts (1), these
 168 filtrations are interleaved, provided that the eigenvalues of the matrices in Σ are positive.

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

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

172 The condition of non-negative weights is not too restrictive since for general weights, it suffices
 173 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.
 174 Then, the condition on t' becomes $\min_{i \in I} \omega_i < t' \leq \frac{\lambda_{\min}}{\lambda_{\max}} \frac{d+1}{2d} t + \left(1 - \frac{\lambda_{\min}}{\lambda_{\max}} \frac{d+1}{2d}\right) \min_{i \in I} \omega_i$.

175 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
 176 optimal. When \mathbf{m} is the set of vertices of a regular d -simplex, the left inclusion is an equality.

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

183 2.3 Examples of filtrations based on robust power functions

184 2.3.1 Isotropic robust power functions

185 Set \mathbb{X} , a set of n points generated on the neighborhood of a compact subset \mathcal{X} of \mathbb{R}^d . In order
 186 to face the non robustness of the distance function to \mathbb{X} , $d_{\mathbb{X}}$, Chazal et al. have introduced
 187 the notion of distance-to-measure (DTM), in [10]. The DTM is a counterpart of $d_{\mathbb{X}}$ robust
 188 to noise and outliers. Its robustness follows from some parameter $k \in \llbracket 1, n \rrbracket$, the number
 189 of nearest-neighbors X^1, X^2, \dots, X^k of x in \mathbb{X} , used to estimate $d_{\mathbb{X}}(x)$. The DTM $d_{\mathbb{X},k}$
 190 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 $m_{x,k} = \sum_{i=1}^k X^i$, the
 191 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.
 192 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
 193 the barycenter of the point cloud \mathbb{X} , up to some factor, which is robust, but very poor in
 194 terms of topological information. The DTM is actually a weighted power function [18]:

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

196 This follows from the fact that the mean distance between x and its k nearest neighbors is
 197 not larger than the mean distance between x and the k nearest neighbors of any other point
 198 $y \in \mathbb{R}^d$. This infimum is actually a minimum over a set of c points $\mathbf{y} = (y_i)_{i \in [1, c]}$ in \mathbb{R}^d , with
 199 c of order $\binom{n}{k}$. A power approximation of the DTM, the k -witnessed distance, was defined
 200 in [18] by replacing \mathbb{R}^d by \mathbb{X} in (4). Its sublevel sets are unions of n balls. An approximation
 201 of the DTM with c (possibly much smaller than n) balls, the c -PDTM, was defined in [6], by
 202 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
 203 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, $\mathbf{y}_{c,k}$
 204 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
 205 graph of the c -PDTM is necessarily above the graph of the DTM. According to [6], for a
 206 sample on a regular d' -dimensional manifold, c can be chosen of order $n^{\frac{d'}{d'+4}}$, which is much
 207 smaller than n . Moreover, the c -PDTM is a good approximation of $d_{\mathcal{X}}^2$, despite noise.

208 2.3.2 An anisotropic robust power function

209 An anisotropic version of the c -PDTM has been introduced in [4], the c -power likelihood to
 210 measure (c -PLM). It consists in replacing Euclidean norms with Mahalanobis norms. For
 211 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
 212 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
 213 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
 214 norm. Set $\boldsymbol{\theta}_{c,k}$, a family of c pairs $(y, \Sigma) \in \mathbb{R}^d \times \mathcal{M}_d$ that minimizes (or which criterion is as
 215 close as possible to the optimal criterion, in case of non existence of a minimum) the following
 216 “k-means”-type criterion $R_{c,k}$ among all $\boldsymbol{\theta}$ s of cardinality c : $R_{c,k}(\boldsymbol{\theta}) = \sum_{i=1}^n \min_{(y,\Sigma) \in \boldsymbol{\theta}} \|X_i -$
 217 $m_{y,\Sigma,k}\|_{\Sigma^{-1}}^2 + v_{y,\Sigma,k} + \log(\det(\Sigma))$. The term $\log(\det(\Sigma))$ prevents optimal covariance matrices
 218 to be degenerated, with Σ^{-1} going to 0. In some sense, minimizing such a criterion boils
 219 down to fit Gaussian distributions to the data set \mathbb{X} , at best. The c -PLM is the power
 220 function defined from $\boldsymbol{\theta}_{c,k}$ by: $x \mapsto \min_{(y,\Sigma) \in \boldsymbol{\theta}_{c,k}} \|x - m_{y,\Sigma,k}\|_{\Sigma^{-1}}^2 + v_{y,\Sigma,k} + \log(\det(\Sigma))$. A
 221 modification of the criterion $R_{c,k}$ has been introduced in [4], to remove some datapoints
 222 ($|\mathbb{X}| - sig$ for some parameter sig), when \mathbb{X} is corrupted with outliers. The criterion is given by
 223 $R_{c,k,sig}(\boldsymbol{\theta}) = \min_{(i_1, i_2, \dots, i_{sig}) \in [1, |\mathbb{X}|]} \sum_{j=1}^{sig} \min_{(y,\Sigma) \in \boldsymbol{\theta}} \|X_{i_j} - m_{y,\Sigma,k}\|_{\Sigma^{-1}}^2 + v_{y,\Sigma,k} + \log(\det(\Sigma))$.
 224 Iterative Lloyd-type algorithms [22] provide local minima $\tilde{\boldsymbol{\theta}}_{c,k}$ and $\tilde{\boldsymbol{\theta}}_{c,k,sig}$ for the criteria
 225 $R_{c,k}$ and $R_{c,k,sig}$ [4]. These algorithms run in $O(ncd^2 + nkd^2 + n \log(n)c)it$ operations, with
 226 it the number of iterations of the algorithm. They consist, given $\boldsymbol{\theta} = (\mathbf{y}, \boldsymbol{\Sigma})$, in splitting the
 227 space \mathbb{R}^d into weighted $\boldsymbol{\Sigma}$ -curved Voronoi cells, replacing centers \mathbf{y} by the centroid of the cells,
 228 and updating the matrices in $\boldsymbol{\Sigma}$ by a close formula from the points in the cells and ellipsoids.
 229 To compute $\tilde{\boldsymbol{\theta}}_{c,k,sig}$, a trimming step is added at each iteration. For clustering, disposing of
 230 a local minimum is enough, as enhanced in the numerical illustration section, since we can
 231 remove bad centers in $\tilde{\boldsymbol{\theta}}_{c,k}$ or in $\tilde{\boldsymbol{\theta}}_{c,k,sig}$ with the parameter *Threshold* in Algorithm 1.

232 3 Persistence-based clustering from power-functions-based filtrations

233 3.1 Persistence for power-functions-based filtrations

234 Set $f_{\mathbf{m},\boldsymbol{\omega},\boldsymbol{\Sigma}} : x \in \mathbb{R}^d \mapsto \min_{i \in I} \|x - m_i\|_{\Sigma_i^{-1}}^2 + \omega_i$, an anisotropic power-function indexed by
 235 a set $I = [1, c]$ and with the ω_i s sorted in non-decreasing order. As above-mentioned, the
 236 sublevel sets $V^t = f_{\mathbf{m},\boldsymbol{\omega},\boldsymbol{\Sigma}}^{-1}((-\infty, t])$ are unions of at most c ellipsoids $\mathcal{E}_i^t = B_{\Sigma_i}(m_i, \sqrt{t - \omega_i})$,
 237 non empty as soon as $t \geq \omega_i$. In particular, each sublevel set of $f_{\mathbf{m},\boldsymbol{\omega},\boldsymbol{\Sigma}}$ contains at most c
 238 connected components. Each connected component of V^t , V_i^t is indexed by the smallest index

239 $i \in I$ such that m_i belongs to the component. With a language abuse, we call connected
 240 component V_i , the family of connected components $(V_i^t)_{t \in T}$ that gets born at time $t = b_i = \omega_i$
 241 and dies at a time $t = d_i$ when V_i^t merges with another connected component V_j^t for some
 242 $j \leq i$. Note that $d_1 = \infty$. The lifetime of the component V_i^t , $d_i - b_i$, is called persistence
 243 or prominence of the component i . This merging information is encoded in a barcode or a
 244 dendrogram. In these two representations, each line is associated to a component V_i , has
 245 length $d_i - b_i$, and begins at the height b_i . The dendrogram is obtained from the barcode by
 246 linking the bars associated to merging components, at a height given by the merging time.

247 When \mathbf{m} is a point set \mathbb{X} , $\Sigma_i = I_d$ and $\omega_i = 0$ for every i , clustering points accordingly
 248 to the connected components of V^t boils down to the classical single-linkage clustering
 249 procedure, with $t > 0$, calibrated in accordance with the dendrogram. This procedure is not
 250 robust to outliers. In this paper, we consider an adjacent procedure, similar to the ToMATo
 251 algorithm [12], based on the prominence of components. To be precise, in the clustering
 252 scheme, a component V_i cannot merge with another component V_j at a time t larger than
 253 $\omega_i + Stop$, for some parameter $Stop$. In other words, components with large prominence will
 254 never die in this clustering procedure. This is the purpose of Algorithm 1 in the next section.

255
 256 In order to better visualize the prominence of the components, we represent their lifetimes
 257 in a persistence diagram. A persistence diagram is a multiset of points $(b_i, d_i) \in \mathbb{R}^2$ that lie
 258 above the diagonal $b = d$. Each point (b_i, d_i) is associated to a connected component V_i . The
 259 notion of persistence diagram was introduced by Edelsbrunner et al. in [16], in the broader
 260 framework of homology, and allows to compute lifetimes of additional features such as loops,
 261 voids etc. It is defined for filtrations that are regular enough, on triangulable spaces such
 262 as \mathbb{R}^d . The proper notion of regularity is the notion of q -tameness [11]. In [7, Proposition
 263 3.5], Buchet et al. proved that the DTM is q -tame. The proof of [7] can be straightforwardly
 264 adjusted for distance functions to compact sets and most importantly, for anisotropic power
 265 functions, provided that the eigenvalues of the matrices Σ_i are all positive.

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

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

275 According to Theorem 3, the persistence diagram of any anisotropic power function
 276 $f_{\mathbf{m}, \omega, \Sigma}$ that is $\epsilon - \|\cdot\|_\infty$ close to $d_{\mathcal{X}}$ is ϵ -bottleneck close to the persistence diagram of the
 277 sublevel sets of $d_{\mathcal{X}}$. Consequently, prominence of the connected components of \mathcal{X} can be
 278 deduced from the diagram associated to $f_{\mathbf{m}, \omega, \Sigma}$, for ϵ small enough. This bottleneck closeness
 279 occurs with large probability for a regular manifold \mathcal{X} for the c -PDTM built from a noisy
 280 sample from \mathcal{X} , according to [6]. No such result has been proved yet for the c -PLM. Anyway,
 281 intuitively, its sublevel sets are good approximations of the manifold \mathcal{X} , with the advantage
 282 that they are made of less ellipsoids, and that these ellipsoids are oriented accordingly to
 283 the manifold, i.e. with large eigenvalues on the tangent space and small eigenvalues on its
 284 orthogonal. This will be confirmed in the numerical illustrations section.

285 By construction, the persistence diagram (for connected components) associated to the
 286 filtration of the sublevel sets of $f_{\mathbf{m}, \omega, \Sigma}$ coincides with the persistence diagram associated to

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

292 3.2 An algorithm for persistence-based clustering

293 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,
 294 we define an algorithm, strongly inspired from the ToMATo algorithm [12]. The clustering
 295 scheme is guided by the persistence of the connected components in $(\mathcal{G}^t)_{t \in \mathbb{R}}$, and preserves
 296 components with large prominence. We assume that the nodes of \mathcal{G} are labeled such that the
 297 node labeled i gets born before the node labeled j , when $i \leq j$. The procedure is as follows.
 298 A connected component gets born when a node gets born, with the same label. A component
 299 changes of label at each time t for which it merges with a component with smaller label in \mathcal{G}^t ,
 300 unless its prominence is larger than some parameter $Stop$. The prominence of a node or a
 301 component is defined as the lifetime of the component in the filtration (i.e. the elapsed time
 302 between the birth of the node and the time t such that a node with smaller index is present
 303 in its connected component in \mathcal{G}^t). The resulting clustering is given by the label of the nodes
 304 at time $t = +\infty$. It contains exactly labels of edges with a prominence larger than $Stop$. In
 305 this clustering scheme, we decide that nodes born after some time parameter $Threshold$ are
 306 not relevant; they are removed. This procedure is implemented in Algorithm 1.

307 **Algorithm 1** Persistence-based Clustering Algorithm

```

307 Data: Mat, Threshold, Stop
308 Result: Color, Birth, Death
309 Initialization ;
310  $c \leftarrow \max\{i \mid \text{Mat}[i,i] \leq \text{Threshold}\}$  ;  $\text{Mat} \leftarrow \text{Mat}[1:c,1:c]$  ;
311  $\text{Birth} \leftarrow [\text{Mat}[i,i] \text{ for } i \text{ in } 1:c]$  ;  $\text{Death} \leftarrow [\infty \text{ for } i \text{ in } 1:c]$  ;
312  $\text{indice} \leftarrow 1$  ;  $I \leftarrow 1$  ;  $\text{time} \leftarrow \text{Mat}[I,I]$  ;  $\text{Color} \leftarrow []$  ;
313 while  $\text{time} < \infty$  do
314   if  $\text{time} = \text{Mat}[I,I]$  then
315     Component  $I$  appears ;
316      $\text{indice} \leftarrow \text{indice} + 1$  ;  $\text{Mat}[I,I] \leftarrow \infty$  ;  $\text{Color}[I] \leftarrow I$ ;
317   else
318      $(\text{col\_max}, \text{col\_min}) \leftarrow (\max(\text{Color}[I], \text{Color}[J]), \min(\text{Color}[I], \text{Color}[J]))$ ;
319     if  $\text{time} - \text{Birth}[\text{col\_max}] \leq \text{Stop}$  then
320       Components  $\text{col\_max}$  and  $\text{col\_min}$  merge ;
321       Replace all entries  $\text{col\_max}$  by  $\text{col\_min}$  in  $\text{Color}$  ;
322        $\text{Death}[\text{col\_max}] \leftarrow \text{time}$  ;
323     else
324       Component  $\text{col\_max}$  will never die ;
325     end
326      $\text{Mat}[i,j] \leftarrow \infty$  for every  $i, j \leq \min(\text{indice}, c)$  such that
327      $(\text{Color}[i], \text{Color}[j]) \in \{(\text{col\_min}, \text{col\_max}), (\text{col\_max}, \text{col\_min})\}$ ;
328   end
329    $I, J \leftarrow \arg \min_{i,j \leq \min(\text{indice}, c)} \text{Mat}[i,j]$  ;  $\text{time} \leftarrow \text{Mat}[I,J]$ 
330 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 ω' and Σ' 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}' , ω' and Σ' : $x \in \mathbb{X}$ has the same label as m'_i if $\|x - m'_i\|_{\Sigma'_i}^2 + \omega'_i \leq \|x - m'_j\|_{\Sigma'_j}^2 + \omega'_j$ for every j .

Since $f_{\mathbf{m}^*, \omega^*, \Sigma^*}$ 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 $f_{\mathbf{m}', \omega', \Sigma'}(x)$ is the largest. Note that a power function is homogeneous to the square of a distance function. Therefore, for positive weights ω , it could be more appropriate to consider the filtration of sublevel sets of $\sqrt{f_{\mathbf{m}, \omega, \Sigma}}$ instead of $f_{\mathbf{m}, \omega, \Sigma}$.

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

354 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 .

355 The DTM-filtration [1] corresponds to the 1-skeleton of the nerve of the union of balls
 356 $(\bigcup_{x \in \mathbb{X}} \bar{B}(x, r_t(x)))_{t > 0}$ 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
 357 $t \geq d_{\mathbb{X},k}(x)$, for some $p \geq 1$ and with the convention that $\bar{B}(x, -\infty)$ is empty. In Table 1, we
 358 give the coefficients for $p = 1$. The DTM-filtration with $p = 2$ was actually introduced in [7],
 359 leading to what we call Power filtration, which coincides with the sublevel-sets filtration of
 360 the square of a power distance. We also consider additional power-functions-based filtrations,
 361 from the k -witnessed distance [18], the c -PDTM [6] and the c -PLM [4].

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

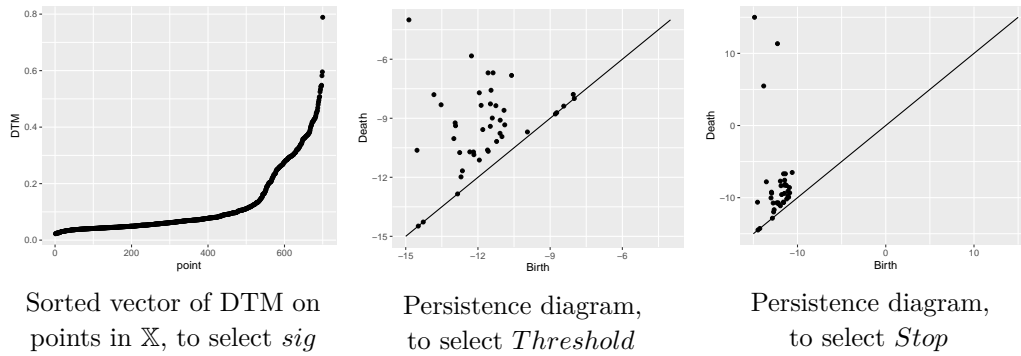
Method	$t_{i,i}$	$t_{i,j}$ for $i < j$
ToMATo	$f(i)$	$\max(f(i), f(j))(\mathbb{1}_{[i,j] \in \mathcal{G}})^{-1}$
DTM-filtration	$f(i)$	$\left(\frac{\ X_i - X_j\ + f(i) + f(j)}{2} \right) \mathbb{1}_{\ X_i - X_j\ > f(i) - f(j) } + f(i) \mathbb{1}_{f(i) - f(j) \geq \ X_i - X_j\ }$
$f_{\mathbf{m},\omega}$	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 $(\mathbf{m}, \omega, \Sigma) = (m_{y,\Sigma,k}, v_{y,\Sigma,k} + \log(\det(\Sigma)), \Sigma)_{(y,\Sigma) \in \theta_{c,k}}$

374 4 Numerical illustrations

375 4.1 A complete illustration of the method

376 Consider the target \mathcal{X} , a set of three curves in \mathbb{R}^2 . We generate $\mathbb{X} = (X_i)_{i \in \llbracket 1, N_s + N_o \rrbracket}$,
 377 a set of $N_s = 500$ signal points $(X_i = Y_i + Z_i)_{i \in \llbracket 1, N_s \rrbracket}$, with Y_i uniform on \mathcal{X} and Z_i
 378 Gaussian with standard deviation $\sigma = 0.02$; corrupted by $N_o = 200$ outliers, uniform on
 379 $[-1.5, 2.5]^2$. We compare the clustering scheme based on Algorithm 1 with the sublevel
 380 sets of the c -PLM, to the target labels in Figure 2 (left). Parameters are set to $c = 50$
 381 centers, $k = 10$ nearest neighbors, $sig = 520$ points to consider as signal, and $it = 100$
 382 iterations and $n_ini = 10$ initializations to compute a suitable local optimum $\tilde{\theta}_{c,k,sig}$ of the
 383 c -PLM-criterion $R_{c,k,sig}$. Since the DTM $d_{\mathbb{X},k}$ is large for outliers, we select sig from the
 384 curve $([d_{\mathbb{X},k}(X_i), i \in \llbracket 1, N_s + N_o \rrbracket])$ in non-decreasing order), as the point of slope break; see
 385 Figure 1 (left). The DTM can be replaced by any not-trimmed approximation of the c -PLM.

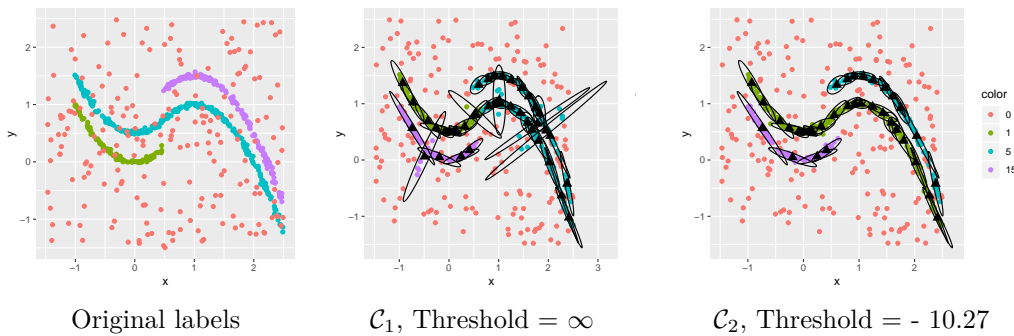
387 We run Algorithm 1 a first time with the parameters $Threshold = \infty$ and $Stop = \infty$, and
 388 display the persistence diagram $(Birth[i], Death[i])_{i \in \llbracket 1, c \rrbracket}$, in Figure 1 (middle). In order to
 389 have 3 clusters, we select $Stop = 5.62$, the height of a line parallel to the diagonal, separating
 390 3 points from the others. We run Algorithm 1 a second time with this new parameter, which
 391 results in the clustering \mathcal{C}_1 of Figure 2 (middle). A sublevel set of the function $f_{\tilde{\theta}_{c,k}}$ is
 392 represented by the union of ellipses. Note that some ellipses have a bad position. This results
 393 in a bad clustering. We use the parameter $Threshold$ to remove them. In Figure 1 (middle),
 394 6 points are on the right side, separated from the other points with a vertical line (of abscissa
 395 -10.27). Then, we run Algorithm 1 with $Threshold = -10.27$ and $Stop = \infty$. According to
 396 the persistence diagram in Figure 1 (right), since 3 points are well-separated from the other



386 **Figure 1** Parameters selection heuristics

397 ones with a large band parallel to the diagonal (containing a line parallel to the diagonal,
 398 with height 12), we recover the number of clusters, 3, and set $Stop = 12$. The clustering \mathcal{C}_2
 399 obtained with $Threshold = -10.27$ and $Stop = 12$ is represented in Figure 2 (right). The
 400 bad ellipses have been removed. Denote by $\tilde{\theta}'_{c,k,sig}$, the subfamily of $\tilde{\theta}_{c,k,sig}$ made of centers
 401 not removed by the procedure. The color of any point x in Figure 2 (right) is given by the
 402 label in $Color$ (label returned by the Algorithm 1) of its associated center (y, Σ) in $\tilde{\theta}'_{c,k,sig}$.
 403 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
 404 labels of the $|\mathbb{X}| - sig$ points with largest $f_{\tilde{\theta}'_{c,k,sig}}$ -value are set to 0.

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

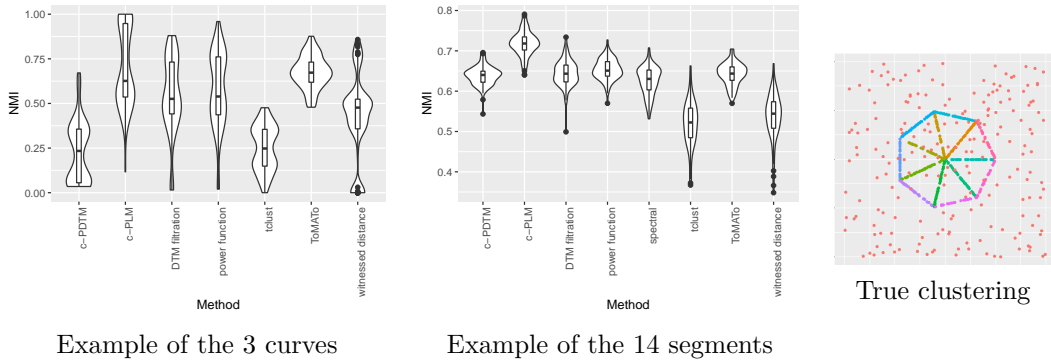


407 **Figure 2** Two resulting clusterings, with ellipses

408 We compare the performance of the two clusterings \mathcal{C}_1 and \mathcal{C}_2 . In terms of outliers
 409 detection, this can be assessed via the proportion of signal points labeled as outliers (0.034
 410 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 ,
 411 0.14 for \mathcal{C}_2). As expected from Figure 2, removing bad ellipses reduces these proportions
 412 and thus improves the outliers detection performance. In terms of clusters recovering, the
 413 normalized mutual information (NMI) is classically used. It equals 1 for a perfect clustering
 414 and 0 for a terrible clustering. When considering outliers as a cluster with label 0, we got
 415 $NMI = 0.586$ for \mathcal{C}_1 and $NMI = 0.841$ for \mathcal{C}_2 . The NMI computed on the signal points
 416 labeled as signal points is $NMI = 0.634$ for \mathcal{C}_1 and $NMI = 1$ for \mathcal{C}_2 , a perfect clustering.

417 **4.2 Comparison of the different methods on synthetic datasets**

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



Example of the 3 curves

Example of the 14 segments

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

426 For the three curves, the parameter r for ToMATo is chosen such that the graph is not
 427 connected, the clusterings are acceptable but have more than 3 clusters. The c -PLM often
 428 performs perfectly, and sometimes performs poorly, since the number of bad ellipses removed
 429 is fixed to 10 and not calibrated according to the heuristics, and there is some instability. We
 430 observe the same clustering problem as in Figure 2 (middle) for the other methods since the
 431 lines are close, compared to the distance between sample points from the same line. For the
 432 polygonal line of 14 segments, all methods except the c -PLM and `tblust` put centers of clusters
 433 on massive parts of \mathbb{X} (the center and the intersections of 3 segments). For the c -PLM and
 434 `tblust`, most clusters coincide with segments. Nonetheless, there is some instability (much
 435 less pronounced for the c -PLM), since the algorithms are based on local minimizers.

436 **4.3 Applications to real datasets**437 **4.3.1 Recovering fleas species, based on 6 measurements**

438 We picked the dataset flea from the R-package `tourr` [29], initially from [23]. This dataset
 439 contains records of 6 measurements for 74 male insects from the Palaearctic, from three
 440 different species : *Heptapotamica*, *Concinna*, *Heikertingeri*. The variables correspond to
 441 measurements on the tarsus, the aedeagus and the head. We normalized data so that the
 442 mean and variance of each of the 6 variables are respectively 0 and 1. In Table 2, we computed
 443 the NMI between the true species and the clustering returned by different methods. We ran
 444 each algorithm 10 times with at most 100 iterations. For every k -nearest-neighbours-based
 445 algorithm, we set $k = 10$. For ToMATo, we set $r = 1.9$ so that the graph is connected ;
 446 for the c -PLM and the c -PDTM, $c = 50$ and for `dbscan`, $eps = 1.5$ and $minPts = 10$. The
 447 3-PDTM and 3-PLM methods consist in clustering data according to the weighted Voronoi
 448 cells given by the optimal centers and covariance matrices.

449 ■ **Table 2** NMI between clustering of fleas and their true specie

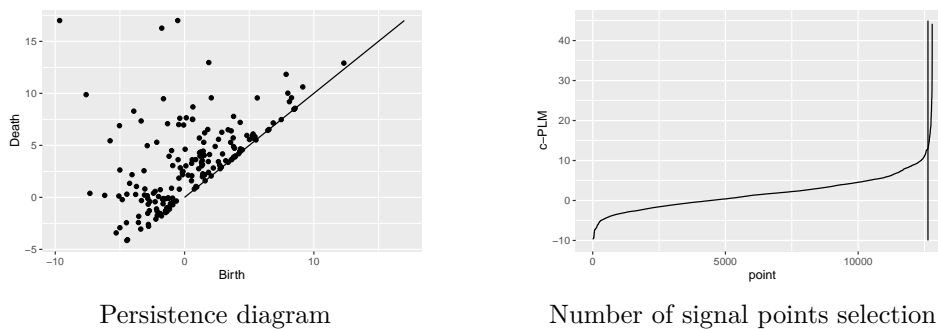
450	Without	<i>k</i> -means	tclust	DBSCAN	Spectral	3-PLM	3-PDTM
451	Algorithm 1	0.825	0.848	0.647	1	1	1
452	With	ToMATo	Witnessed	power	DTM-filt.	<i>c</i> -PLM hier.	<i>c</i> -PDTM hier.
453	Algorithm 1	0.628	0.906	1	1	1	1

454 The methods based on the decomposition of \mathbb{R}^6 into 3 (weighted and/or curved) Voronoi
 455 cells are efficient: at most 3 bad labels for *k*-means and tclust and all labels correct for their
 456 “robust” versions, the 3-PDTM and the 3-PLM. The perfect performance of these two last
 457 functions is due to the weights that force the centers of cells to lie in massive areas for \mathbb{X} . The
 458 bad performance of ToMATo is due to the difficulty to select the parameter *r* for the Rips
 459 graph, the small number of points, and the fact that the inverse of the DTM should be used
 460 instead of the DTM, as recommended by the authors. Nonetheless, we made the choice to
 461 use the DTM since the other methods (witnessed distance, power function, DTM-filtration,
 462 *c*-PLM and *c*-PDTM) are based on filtrations from approximations of the DTM, and almost
 463 all of these methods perform perfectly. The method dbscan performs poorly since it labels
 464 14 points as outliers. Nonetheless, the points considered as signal are well clustered.

465 **4.3.2 Clustering a earthquake dataset**

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

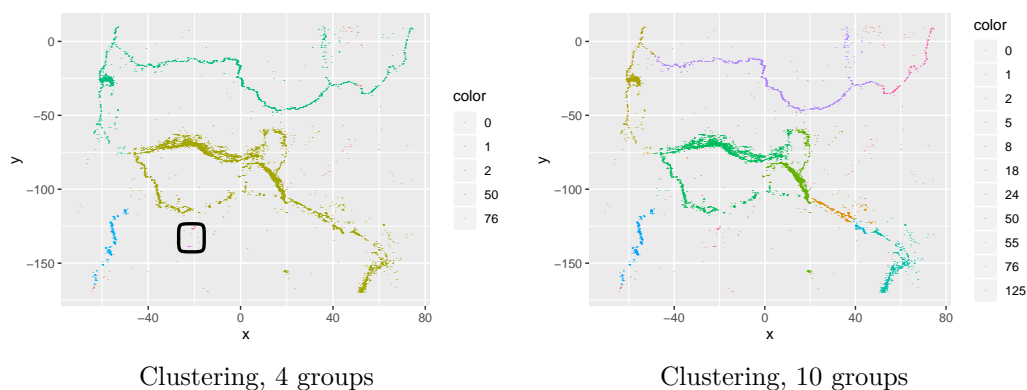
469 We used Algorithm 1 with an approximation of the *c*-PLM based on a sub-sample of
 470 2000 points from the dataset, with parameters *c* = 200, *k* = 10 and for *it* = 50 iterations.
 471 We restricted matrices Σ to have eigenvalues smaller than 50 by thresholding them. The
 472 persistence diagram in Figure 4 suggests that the dataset has 4 or 10 clusters. Moreover, the
 473 curve of the sorted values of the *c*-PLM approximation on the pointset in Figure 4 suggests
 474 to keep *sig* = 12250 points as signal points. See Figure 5 for the corresponding clustering.



475 ■ **Figure 4** Parameters selection heuristics.

477 **References**

478 1 Hirokazu Anai, Fr d ric Chazal, Marc Glisse, Yuichi Ike, Hiroya Inakoshi, Rapha l Tinarrage,
 479 and Yuhei Umeda. DTM-based filtrations. In *35th International Symposium on Computational*



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

- 480 *Geometry*, volume 129 of *LIPICs. Leibniz Int. Proc. Inform.*, pages Art. No. 58, 15. Schloss
 481 Dagstuhl. Leibniz-Zent. Inform., Wadern, 2019.
- 482 2 Arindam Banerjee, Srujana Merugu, Inderjit S. Dhillon, and Joydeep Ghosh. Clustering
 483 with bregman divergences. *J. Mach. Learn. Res.*, 6:1705–1749, December 2005. URL: <http://dl.acm.org/citation.cfm?id=1046920.1194902>.
- 484 3 Gregory Bell, Austin Lawson, Joshua Martin, James Rudzinski, and Clifford Smyth. Weighted
 485 persistent homology. *Involve*, 12(5):823–837, 2019. URL: <https://doi.org/10.2140/involve.2019.12.823>, doi:10.2140/involve.2019.12.823.
- 486 4 Claire Br echeteau. Robust shape inference from a sparse approximation of the gaussian
 487 trimmed loglikelihood. Unpublished, 2018.
- 488 5 Claire Br echeteau, Aur elie Fischer, and Cl ement Levrard. Robust bregman clustering. In
 489 revision, 2018.
- 490 6 Claire Br echeteau and Cl ement Levrard. A k-points-based distance for robust geometric
 491 inference. To appear in *Bernoulli*, 2017.
- 492 7 Micka el Buchet, Fr ed eric Chazal, Steve Y. Oudot, and Donald R. Sheehy. Efficient and
 493 robust persistent homology for measures. *Comput. Geom.*, 58:70–96, 2016. URL: <https://doi.org/10.1016/j.comgeo.2016.07.001>, doi:10.1016/j.comgeo.2016.07.001.
- 494 8 Micka el Buchet, Tamal K. Dey, Jiayuan Wang, and Yusu Wang. Declutter and resample:
 495 towards parameter free denoising. *J. Comput. Geom.*, 9(2):21–46, 2018.
- 496 9 Fr ed eric Chazal, David Cohen-Steiner, Marc Glisse, Leonidas J. Guibas, and Steve Y. Oudot.
 497 Proximity of persistence modules and their diagrams. In *Proceedings of the Twenty-fifth
 498 Annual Symposium on Computational Geometry, SCG ’09*, pages 237–246, New York, NY,
 499 USA, 2009. ACM. URL: <http://doi.acm.org/10.1145/1542362.1542407>, doi:10.1145/
 500 1542362.1542407.
- 501 10 Fr ed eric Chazal, David Cohen-Steiner, and Quentin M erigot. Geometric Inference for Measures
 502 based on Distance Functions. *Foundations of Computational Mathematics*, 11(6):733–751,
 503 2011. URL: <https://hal.inria.fr/inria-00383685>, doi:10.1007/s10208-011-9098-0.
- 504 11 Fr ed eric Chazal, Vin de Silva, Marc Glisse, and Steve Oudot. *The structure and stability of
 505 persistence modules*. SpringerBriefs in Mathematics. Springer, [Cham], 2016. URL: <https://doi.org/10.1007/978-3-319-42545-0>, doi:10.1007/978-3-319-42545-0.
- 506 12 Fr ed eric Chazal, Leonidas J. Guibas, Steve Y. Oudot, and Primoz Skraba. Persistence-
 507 based clustering in Riemannian manifolds. *J. ACM*, 60(6):Art. 41, 38, 2013. URL: <https://doi.org/10.1145/2535927>, doi:10.1145/2535927.
- 508 13 David Cohen-Steiner, Herbert Edelsbrunner, and John Harer. Stability of persistence di-
 509 agrams. *Discrete Comput. Geom.*, 37(1):103–120, 2007. URL: <https://doi.org/10.1007/s00454-006-1276-5>, doi:10.1007/s00454-006-1276-5.

- 516 14 J. A. Cuesta-Albertos, A. Gordaliza, and C. Matrán. Trimmed k -means: an attempt to
517 robustify quantizers. *Ann. Statist.*, 25(2):553–576, 1997. URL: [https://doi.org/10.1214/
518 aos/1031833664](https://doi.org/10.1214/aos/1031833664), doi:10.1214/aos/1031833664.
- 519 15 Vin de Silva and Robert Ghrist. Coverage in sensor networks via persistent homology.
520 *Algebr. Geom. Topol.*, 7:339–358, 2007. URL: <https://doi.org/10.2140/agt.2007.7.339>,
521 doi:10.2140/agt.2007.7.339.
- 522 16 Herbert Edelsbrunner, David Letscher, and Afra Zomorodian. Topological persistence and sim-
523 plification. *Discrete Comput. Geom.*, 28(4):511–533, 2002. Discrete and computational
524 geometry and graph drawing (Columbia, SC, 2001). URL: [https://doi.org/10.1007/
525 s00454-002-2885-2](https://doi.org/10.1007/s00454-002-2885-2), doi:10.1007/s00454-002-2885-2.
- 526 17 Heinrich Fritz, Luis A. Garcia-Escudero, and Agustin Mayo-Iscar. tclust: An R package for
527 a trimming approach to cluster analysis. *Journal of Statistical Software*, 47(12):1–26, 2012.
528 URL: <http://www.jstatsoft.org/v47/i12/>.
- 529 18 Leonidas Guibas, Dmitriy Morozov, and Quentin Mérigot. Witnessed k -distance. *Discrete*
530 *Comput. Geom.*, 49(1):22–45, 2013. URL: <https://doi.org/10.1007/s00454-012-9465-x>,
531 doi:10.1007/s00454-012-9465-x.
- 532 19 Michael Hahsler, Matthew Piekenbrock, and Derek Doran. dbscan: Fast density-based
533 clustering with R. *Journal of Statistical Software*, 91(1):1–30, 2019. doi:10.18637/jss.v091.
534 i01.
- 535 20 Allen Hatcher. *Algebraic topology*. Cambridge University Press, Cambridge, 2002.
- 536 21 Alexandros Karatzoglou, Alex Smola, Kurt Hornik, and Achim Zeileis. kernlab – an S4
537 package for kernel methods in R. *Journal of Statistical Software*, 11(9):1–20, 2004. URL:
538 <http://www.jstatsoft.org/v11/i09/>.
- 539 22 Stuart P. Lloyd. Least squares quantization in PCM. *IEEE Trans. Inform. Theory*, 28(2):129–
540 137, 1982. URL: <https://doi.org/10.1109/TIT.1982.1056489>, doi:10.1109/TIT.1982.
541 1056489.
- 542 23 Alexander A. Lubischew. On the use of discriminant functions in taxonomy. *Biometrics*, pages
543 455–477, 1962.
- 544 24 J. MacQueen. Some methods for classification and analysis of multivariate observations.
545 In *Proceedings of the Fifth Berkeley Symposium on Mathematical Statistics and Probability,*
546 *Volume 1: Statistics*, pages 281–297, Berkeley, Calif., 1967. University of California Press.
547 URL: <https://projecteuclid.org/euclid.bsm/1200512992>.
- 548 25 Stephen B Pope. Algorithms for ellipsoids. Technical Report FDA-08-01, Sibley School of
549 Mechanical & Aerospace Engineering, Cornell University Ithaca, New York 14853, 2008.
- 550 26 P. J. Rousseeuw and A. M. Leroy. *Robust Regression and Outlier Detection*. John Wiley &
551 Sons, New York, 1987.
- 552 27 Ulrike von Luxburg. A tutorial on spectral clustering. *Stat. Comput.*, 17(4):395–416, 2007.
553 URL: <https://doi.org/10.1007/s11222-007-9033-z>, doi:10.1007/s11222-007-9033-z.
- 554 28 Wenping Wang, Jiaye Wang, and Myung-Soo Kim. An algebraic condition for the separation
555 of two ellipsoids. *Comput. Aided Geom. Design*, 18(6):531–539, 2001. URL: [https://doi.
556 org/10.1016/S0167-8396\(01\)00049-8](https://doi.org/10.1016/S0167-8396(01)00049-8), doi:10.1016/S0167-8396(01)00049-8.
- 557 29 Hadley Wickham, Dianne Cook, Heike Hofmann, and Andreas Buja. tourr: An R package for
558 exploring multivariate data with projections. *Journal of Statistical Software*, 40(2):1–18, 2011.
559 URL: <http://www.jstatsoft.org/v40/i02/>.