

Preliminary results on classification based on Topological Data Analysis(1)

Rolando Kindelan

Computer Science Department, Faculty of Mathematical and Physical Sciences, University of Chile, 851 Beauchef Av. Santiago de Chile, Chile

Center of Medical Biophysics, University of Oriente, Santiago de Cuba, Cuba
rolan2kn@gmail.com

Mauricio Cerda

Laboratory for Scientific Image Analysis (SCIAN-Lab), Biomedical Neuroscience Institute (BNI), Faculty of Medicine, University of Chile, 1027 Independencia Av. Santiago, Chile.

mauriciocerda@med.uchile.cl

Nancy Hitschfeld

Computer Science Department, Faculty of Mathematical and Physical Sciences, University of Chile, 851 Beauchef Av. Santiago, Chile.

nancy@dcc.uchile.cl

1 — Abstract —

2 In recent years, algebraic topology techniques have been applied to data analysis, giving rise to an
3 emerging research field called Topological Data Analysis (TDA). TDA tools are useful for inferring
4 the topology underlying a dataset. On the other hand, it is known that most Machine Learning (ML)
5 techniques cannot capture topological information. Consequently, TDA has been used to create
6 filters and topological descriptors to improve the results of ML methods. However, and despite the
7 good results, these hybrid methods do not to fully exploit the potentials of TDA and hence implicit
8 duplications of computations occur. This paper proposes a TDA-based method to solve a 3-class
9 classification problem. The method constructs a filtered simplicial complex K from a training set S
10 and testing set X . This construction produces a large size $O(2^{|S \cup X|})$ searching space to classify each
11 $x \in X$. We use persistent homology to downsize the searching space by detecting persistent intervals
12 of desired topological features. In consequence, the searching space is reduced to $O(2^{q+1} \cdot |X|)$, with
13 $q \ll |S \cup X|$ the constant dimension of the sub-complex $K_i \subset K$ resulting from selected persistent
14 interval. In the last stage, the algorithm labels each $x \in X$, by using the label-contribution of
15 every 0-simplex $\in K_i$ which shares at least a simplex with x . We develop an experimental proof of
16 concept by comparing our TDA-based classifier with a k-NN classifier on the Iris dataset, using 10
17 evaluation metrics in a repeated-cross validation process. Our method gets 96% accuracy versus
18 97% for k-NN. Our work shows that it is possible to classify only with TDA, with no additional
19 machine learning algorithms. For future work we plan to apply the proposed method to datasets
20 with highly dimensional and noisy samples.

2012 ACM Subject Classification Mathematics of Computing → Point-set Topology; Mathematics of Computing → Algebraic Topology; Theory of computation → Computational geometry; Theory of computation → Machine learning theory; Theory of computation → Semi-supervised learning

Keywords and phrases TDA, TDA-based classifier, simplicial complexes, persistent homology, machine learning

Funding *Rolando Kindelan*: CONICYT 2018/BECA DOCTORADO NACIONAL-21181978

Lines 0

(Submitted to 36th SoCG - Zürich, Switzerland - June 23-26, 2020)



© Rolando Kindelan and Mauricio Cerda and Nancy Hitschfeld;
licensed under Creative Commons License CC-BY

42nd Conference on Very Important Topics (CVIT 2016).

Editors: John Q. Open and Joan R. Access; Article No. 23; pp. 23:1–23:17

Leibniz International Proceedings in Informatics



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

21 **1 Introduction**

22 The processing and extraction of information from large and noisy data sets is a challenging
23 problem in Computer Science. The techniques of algebraic topology have gained the attention
24 of scientists for years, giving rise to an emerging research field called Topological Data Analysis
25 (TDA) [24, 5, 8, 9, 30, 13, 12]. TDA is an approach to infer the topology underlying a
26 dataset by using combinatorial algebraic structures known as simplicial complexes. TDA
27 also involves the computation of invariant properties from continuous transformations of
28 these simplicial complexes: a process known as persistent homology [8, 9, 30, 13].

29 Over several decades the high-dimensionality of datasets and the combinatorial and
30 continuous character of Topology have been issues, making computing persistent homology a
31 challenge. Regarding persistent homology, Edelsbrunner et al. [8] present an efficient algorithm
32 and its visualization as a persistence diagram [8, 30]. Carlson et al. [5] strengthened the
33 mathematical foundations and also proposed another visualization tool called persistence
34 Barcodes [13, 5]. Further developments of the TDA field are derived from those initial works.

35 The construction and representation of simplicial complexes also represent a challenge,
36 as a consequence of their combinatorial nature. Many works have dealt with efficient
37 construction, representation and filtration of simplicial complexes. As a consequence, efficient
38 data structures and algorithms were developed [7, 1, 29, 2, 3, 17, 18], mainly focusing on
39 efficient construction of Čech, Rips and other kinds of simplicial complexes not used here.
40 Theoretical and practical results have been organized in the form of TDA libraries like
41 GUDHI [3, 19], Dionysus, Ripser, DIPA, Perseus and JavaPlex. A complete benchmark of
42 those libraries can be found in Otter et al. [21].

43 A TDA-based method was used in [11] for classifying high-resolution diabetic retinopathy
44 images. They use a preprocessing stage for computing persistent homology to detect
45 topological features encoded into persistence diagrams. A support vector machine (SVM)
46 was used to classify the images according to the persistence descriptors which was used
47 to discriminate between diabetic and healthy patients. They recommend exploring their
48 TDA+SVM method further in larger datasets of high-resolution images.

49 Also TDA has been applied to time-series analyses [6]. One common pipeline is to consider
50 the time-series as a dynamic system and compute the attractors or time-variant of the signal
51 which creates a manifold around the attractors, turning the signal into phase-domain [28,
52 27]. Persistent Homology or another TDA-tool is applied on these phase-space manifold to
53 create topological descriptors [23] and as a final step, a machine learning method is applied
54 like k-NN, CNN, and also SVM. There are other applications of TDA presented in [6] like
55 scientific visualization, bioinformatics, atmospheric and climate data analysis, cosmology
56 and combustion simulations. Furthermore, TDA has been applied in neurosciences [25] and
57 recently in human motion understanding [16, 14, 28].

58 All those examples of TDA-applications have one thing in common: TDA has been used
59 as a preprocessing stage of conventional Machine Learning (ML) algorithms [26]. However,
60 to the best of our knowledge, there are no references to using TDA directly as a classification
61 method.

62 In this paper, we propose a classification method that takes advantage of the TDA
63 information and topological properties. We compare our method to k-NN as a baseline, one
64 of the most used supervised classification methods. This document has been organized as
65 follows. Section 2 exposes the mathematical foundations that we use in this work. Section 3
66 explains the concepts, algorithms, and methodology of our proposed classification method.
67 Next, Section 4 presents our experimental design, evaluation criteria, selected metrics, and

68 our results. Section 5 is where we explain our preliminary results and development decisions
 69 of our method. Section 7 presents several research lines that have arisen from this work. We
 70 present the main conclusions of our work in the Section 6.

71 **2 Mathematical foundations**

72 In this section, we introduce the mathematical definition of simplices, simplicial complex, the
 73 Čech and Rips complexes, the star, and link concepts. We also define persistent homology,
 74 filtration, sub-complex, and filtration levels.

75 **2.1 Simplicial Complexes**

76 Simplicial complexes are combinatorial and algebraic objects which represent a discrete space
 77 homotopically equivalent to data space. There are some related concepts to understand. In a
 78 nutshell, a q -simplex is the convex hull of $q + 1$ affinely independent points $\{s_0, s_1, \dots, s_q\} \in$
 79 $S, S \subset \mathbb{R}^n$. A q -simplex σ has dimension $\dim(\sigma) = q$ and cardinality $|\sigma| = \text{card}(\sigma) =$
 80 $\dim(\sigma) + 1$. A simplex τ defined by $S' \subseteq S$ is a *face* of σ and has σ as a *coface*. A q -simplex
 81 has $\binom{q+1}{d+1}$ faces of dimension d and $\sum_{d=-1}^q \binom{q+1}{d+1} = 2^{q+1}$ faces in total. We symbolize the
 82 *face* and *coface* relationships with $\sigma \geq \tau$ and $\tau \leq \sigma$. So, a simplicial complex K is a finite
 83 collection of simplices such that:

- 84 ■ $\sigma \in K$ and $\tau \leq \sigma \implies \tau \in K$.
- 85 ■ $\sigma_1, \sigma_2 \in K \implies \sigma_1 \cap \sigma_2 \leq \sigma_1, \sigma_2$.

86 The dimension of K is $\dim(K) = \max\{\dim(\sigma) \mid \sigma \in K\}$.

87 There are many known simplicial complexes, though the most popular are the Čech and
 88 Vietoris-Rips complexes. The Čech complex is built by all non-empty intersections of closed
 89 balls $\mathbb{B}_s(\varepsilon)$, with ε radius and centered on each point s from the dataset [13, 8]:

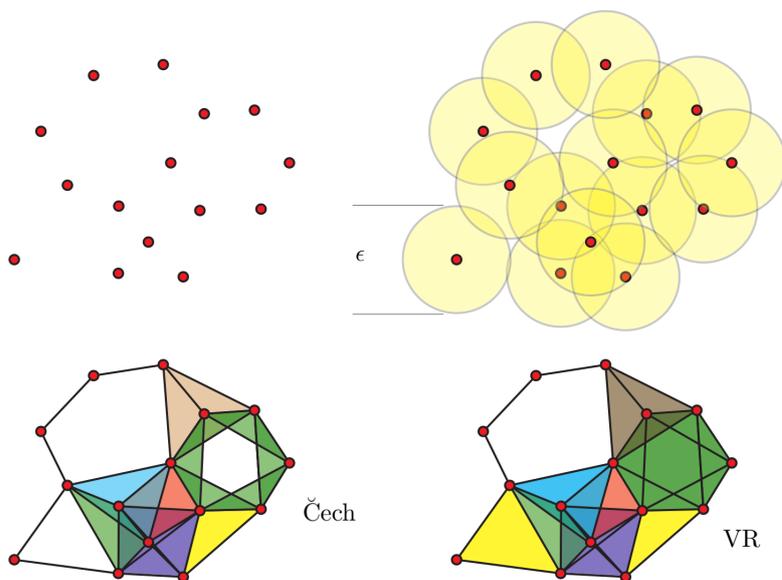
$$90 \quad \check{C}ech(\varepsilon) = \{\sigma \subseteq S \mid \bigcap_{\tau \in \sigma} \mathbb{B}_\tau(\varepsilon) \neq \emptyset\}.$$

91 The Vietoris-Rips (VR) complex from a point set S and ε value is built with all subsets, for
 92 which each minimum enclosing ball has a diameter up to 2ε [8]:

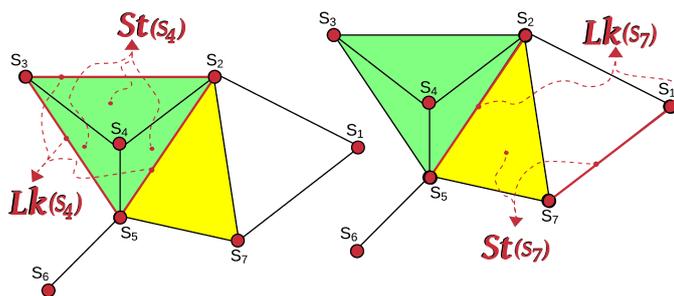
$$93 \quad VR(\varepsilon) = \{\sigma \subseteq S \mid \text{diam}(\sigma) \leq 2\varepsilon\}.$$

94 This implies $\check{C}ech(\varepsilon) \subseteq VR(\varepsilon) \subseteq \check{C}ech(\sqrt{2}\varepsilon)$ a proof is given in [8], this relationship is shown
 95 in Figure 1. The Čech complex is intrinsically a high dimensional simplicial complex. From
 96 a computational sense, VR complex is more feasible (i.e. lower storage and time complexity)
 97 than Čech, even when VR complex has more simplices in general. Compared to Čech, VR
 98 complex does not need to be completely stored, it can be stored like a graph and reconstituted
 99 combinatorially [13].

104 ► **Definition 1 (Star and Link).** Let K be a simplicial complex, and $\tau \in K$ a q -simplex.
 105 The star of τ defined by $St(\tau) = \{\sigma \in K \mid \tau \leq \sigma\}$ is the set of all cofaces of τ [18, 8]
 106 (see Figure 2). The $St(\tau)$ is not a simplicial complex because of the missing faces. If those
 107 faces are added to $St(\tau)$, we get the closed star of τ denoted by $\overline{St}(\tau)$, which is the smallest
 108 simplicial complex that contains the star. The link of τ is a set of simplices in the closed
 109 star that does not share any face with τ , $Lk(\tau) = \{\nu \in \overline{St}(\tau) \mid \nu \cap \tau = \emptyset\}$ [18, 8]. If τ is a
 110 0 -simplex, then $Lk(\tau) = \overline{St}(\tau) - St(\tau)$ (see Figure 2).



94 **Figure 1** From a point set [upper left] a proximity parameter *varepsilon* is applied [upper right]
 95 and two complexes were built: a Čech complex [lower left] and a VR complex [lower right]. There
 96 are differences between $VR(\epsilon)$ and $\check{C}ech(\frac{\epsilon}{2})$. VR has more simplices than Čech as expected, note
 97 all colored 2-simplices. This picture was taken from [13].



118 **Figure 2** Example of $Lk(s_4), St(s_4)$ and $Lk(s_7), St(s_7)$ on a given simplicial complex K .

111 In Figure 2 we present two examples of the St and Lk from points s_4, s_7 from a point set S .
 112 If we expand the solutions we get the following:

113
$$St(s_4) = \{\{s_4\}, \{s_4, s_2\}, \{s_5, s_4\}, \{s_4, s_3\}, \{s_3, s_4, s_2\}, \{s_3, s_5, s_4\}, \{s_4, s_5, s_2\},$$

114
$$\{s_3, s_5, s_2, s_4\}\},$$

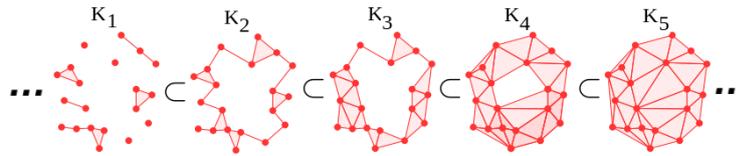
115
$$Lk(s_4) = \{\{s_3\}, \{s_2\}, \{s_5\}, \{s_2, s_3\}, \{s_3, s_5\}, \{s_5, s_2\}, \{s_3, s_5, s_2\}\},$$

116
$$St(s_7) = \{\{s_7\}, \{s_7, s_1\}, \{s_5, s_7, s_2\}\},$$

117
$$Lk(s_7) = \{\{s_1\}, \{s_2\}, \{s_5\}, \{s_5, s_2\}\}.$$

119 2.2 Persistent Homology

120 Persistent homology is a tool to find topological features in a metric space [8, 13, 5]. As
 121 a general rule, the objective of persistent homology is to track how topological features



129 **Figure 3** A fragment example of a simplicial complex filtration..

122 on a topological space appear and disappear when a scale value (usually a radius) varies
 123 incrementally, a process known as filtration [9, 29, 30].

124 **Definition 2 (Sub-complex).** *Let K be a simplicial complex. K' is a sub-complex of K if*
 125 *$K' \subseteq K$ and besides K' is by itself a simplicial complex.*

126 **Definition 3 (Filtration).** *Let K be a simplicial complex. A filtration f is a succession of*
 127 *increasing sub-complexes of K :*

$$128 \quad \emptyset = K_0 \subseteq K_1 \subseteq K_2 \subseteq K_3 \subseteq \dots \subseteq K_n = K.$$

130 *We can understand a filtration as a method to build the whole simplicial complex K from*
 131 *a “family” of sub-complexes incrementally sorted according to some criteria, where each*
 132 *level i corresponds to the “birth” or “death” of a q -simplex set as described in the following*
 133 *definition.*

134 **Definition 4 (Birth and Death).** *birth is a metaphorical concept to describe the filtration*
 135 *level when a set of simplices are created. Similarly, death refers to the filtration level when a*
 136 *set of simplices disappeared. Thus, a persistence interval (birth, death) is the “lifetime” of a*
 137 *given set of simplices Q [8, 9, 5, 29, 30]. Note Q is not a sub-complex because of missing*
 138 *faces with other intervals.*

139 **3 Proposed Classification Method**

140 Let P be a metric space, with every $p \in P$ a data feature value vector in R^n . We split the
 141 dataset P in two sets S, X in which $S \neq \emptyset; X \neq \emptyset; S \cap X = \emptyset$ and $S \cup X = P$. Let S be
 142 the training set and X the testing set to be classified. Let L be a label set. The 2-tuple set
 143 $T = \{(p, l) \mid p \in P; l \in L\}$ relates each element of P with its associated label from L . The
 144 incomplete association set T' is a 2-tuple set with $T' = \{(p, l) \mid (p \in S; l \in L) \vee (p \in X; l = \emptyset)\}$
 145 where each element of S has an associated label from L , but each element of X does not have
 146 any label.

147 **3.1 Definitions**

148 In order to understand and clarify the proposed method, we define several mathematical
 149 concepts.

150 **Definition 5 (Useful-simplex and Non-useful-simplex).** *Let K be a simplicial complex built*
 151 *from $S \cup X$. A q -simplex $\sigma \in K$. Let α, β be two sets where $\alpha \subset S, \beta \subset X, \alpha \subseteq \sigma, \beta \subseteq \sigma$ and*
 152 *$\alpha \cup \beta = \sigma, \alpha \cap \beta = \emptyset$. We say σ is a useful-simplex if $|\alpha| > |\beta|$. In another case, if $|\alpha| \leq |\beta|$,*
 153 *then σ is a non-useful-simplex.*

23:6 Classification based on Topological Data Analysis

154 ► **Definition 6** (Association function, Ψ). Let K be a simplicial complex, σ a q -simplex $\sigma \in K$.
 155 Let $\mathcal{P}(L)$ be the power-set of L . Let \mathbb{T} be a 2-tuple set of associations $\mathbb{T} \in \{T, T'\}$. We define
 156 the association function $\Psi : K \times K \rightarrow \mathcal{P}(L)$ such as:

$$157 \quad \Psi(\sigma, K) = \begin{cases} l \in \mathcal{P}(L) \vee (\sigma, l) \in \mathbb{T} & \text{iff } \text{card}(\sigma) = 1; l \neq \emptyset, \\ \bigcup_{\alpha \in \text{Lk}(\sigma)} \Psi(\alpha, K) & \text{iff } \text{card}(\sigma) = 1; l = \emptyset, \\ \bigcup_{\tau \in \sigma} \Psi(\tau, K) & \text{iff } \text{card}(\sigma) > 1 \end{cases} \quad (1)$$

158 When $\text{card}(\sigma) = 1$, σ is a 0-simplex with, $\sigma \in S$, or (exclusive or) $\sigma \in X$. In the first case it
 159 is obvious that $l \neq \emptyset$ because the labels are always known for all elements in S . In the second
 160 case, $l \neq \emptyset$ if l was previously computed or $\mathbb{T} = T$. In any other case, as $l = \emptyset$, we found l
 161 value by using the union of all associations of each element of $\text{Lk}(\sigma)$.

162 By definition 1, for a simplicial complex K , the intersection of Lk of any element of a
 163 simplex $\sigma \in K$, is necessarily non-empty, therefore all solutions of $\Psi(\sigma, K)$ are overlapping
 164 and it is possible to optimize its computation by using dynamic programming.

165 ► **Definition 7** (Indicator function, \mathcal{I}). Let c be a logical statement. The indicator function \mathcal{I}
 166 is defined by:

$$167 \quad \mathcal{I}(c) = \begin{cases} \text{if } c = \text{True} & 1, \\ \text{if } c = \text{False} & 0. \end{cases} \quad (2)$$

168 Function $\Gamma(\sigma, K)$ returns a vector $V \in \mathbb{R}^{|L|}$ for which each element $v_i \in V$ represents the
 169 amount of apparitions (votes) achieved by label $l_i \in L$ during $\Psi(\sigma, K)$ computation. According
 170 to definition 6, we can also make optimizations on computing $\Gamma(\sigma, K)$ by using dynamic
 171 programming.

172 ► **Definition 8** (Labeling function, Υ). Let K be a simplicial complex, and σ a q -simplex
 173 $\sigma \in K$. Let $\mathcal{M} : \mathbb{R}^{|L|} \rightarrow \mathbb{N} \cup \{0\}$ a function, which by taking a vector $V \in \mathbb{R}^{|L|}$ returns an
 174 integer $0 \leq i < |L|$ where i is the position of maximum value in V . We define the labeling
 175 function of a q -simplex as follows:

$$176 \quad \Upsilon(\sigma, K) = \mathcal{G}(i); i = \mathcal{M}(\Gamma(\sigma, K)). \quad (3)$$

177 Function $\Upsilon(\sigma, K)$ assigns to σ the most voting label $l \in L$ during computation of $\Psi(\sigma, K)$.

178 ► **Definition 9** (Real label list, Y). Let assume a given lexicographic order in X and another
 179 one in L . Y is a list of labels assigned to each element of X defined by:

$$180 \quad Y = (l_i \mid (x_i, l_i) \in T; x_i \in X; l_i \in L).$$

181 we call Y the real label list of X .

182 ► **Definition 10** (Predicted label list, \hat{Y}). Let K be a simplicial complex, T an association
 183 set, T' an uncompleted association set and L a label set. Assuming a lexicographic order in
 184 X , we define the predicted label list \hat{Y} as:

$$185 \quad \hat{Y} = (\hat{y}_i \mid \hat{y}_i = \Upsilon(x_i, K); x_i \in X).$$

186 Note, when we assume $\mathbb{T} = T$ as known associations in basic cases of computation of Ψ , then
 187 $\hat{Y} = Y$, with Y the real labels list from definition 9. In another case, when using $\mathbb{T} = T'$ as
 188 associations then \hat{Y} is a prediction of labels of any value in X .

189 ► **Lemma 11.** *Let P be a point set and S a non-empty subset of P . Let C be a simplicial*
 190 *complex built from P . If K is a simplicial complex built on S with the same construction*
 191 *rules of C , then $K \subseteq C$. In other words: K is a sub-complex of C .*

192 **Proof.** Suppose by contradiction $K \not\subseteq C$. According to section 2.1: we have $\forall \sigma \in K; \tau \in$
 193 $\sigma \implies \tau \in K$. Now, if C is formed by all points in P , then C is a set of subsets of P . Since
 194 C and K have the same formation rules, moreover $S \subseteq P$, then $\forall \sigma \in K \implies \sigma \in C$. As
 195 well as $\forall \sigma \in K; \sigma \in C$ also $K \in C$, which implies $K \subseteq C$. That makes K in a sub-complex
 196 of C . But this contradicts the premise of $K \not\subseteq C$. ◀

197 ► **Definition 12** (Simplicial complexes union operator, \uplus). *Let S, X be two point sets. Let*
 198 *K, A be two simplicial complexes where K was built from S and A built from X respectively.*
 199 *The union operator \uplus is defined by $K \uplus A = C$ where C is a simplicial complex built from*
 200 *$\{S \cup X\}$.*

201 ► **Definition 13** (Filtration level access function, ψ). *Let \mathcal{F} be the set of all possible filtrations*
 202 *on a point set S . Let K be a simplicial complex built from S . The level access function*
 203 *$\psi : \mathcal{F} \times \mathbb{N} \rightarrow K$ is defined by $\psi(f, i) = K_i$. In other words, for a filtration $f \in \mathcal{F}$ and a level*
 204 *i , ψ returns the sub-complex in the i^{th} level of f . If $i \geq n$ with n the maximum level of f ,*
 205 *then $\psi(f, n) = K_n$ is invoked.*

206 Let C be a simplicial complex built from P , let f' be a filtration of C (see Figure 4) with
 207 $n + 1$ levels:

$$208 \quad \emptyset = C_0 \subseteq C_1 \subseteq C_2 \subseteq \dots \subseteq C_n = C,$$

209 and let K be a simplicial complex built from S , with $S \subseteq P$, and let f be a filtration of K
 210 (see Figure 3) with $m + 1$ levels and $m \leq n$:

$$211 \quad \emptyset = K_0 \subseteq K_1 \subseteq K_2 \subseteq \dots \subseteq K_m = K.$$

212 Thus, for each level value i , we can get $C_i = \psi(f', i)$ and $K_i = \psi(f, i)$ according to
 213 Definition 13. As we said in Lemma 11 and definitions 1, 3 and 12 we can affirm $\forall i \leq n$:

$$214 \quad K_i \subseteq C_i, \tag{4}$$

$$215 \quad C_i - \{St(x, C_i) \mid x \in C_0 \wedge x \in X\} = K_i, \tag{5}$$

$$216 \quad K_i \uplus \{St(x, C_i) \mid x \in C_0 \wedge x \in X\} = C_i. \tag{6}$$

217 That is, if we remove from C_i those simplices which contain elements of X , then we can get
 218 K_i . At the same time, by using the union operator of definition 12 and equation 6, if we
 219 include all elements of X in K_i step by step, we can get some sub-complex of C_i .

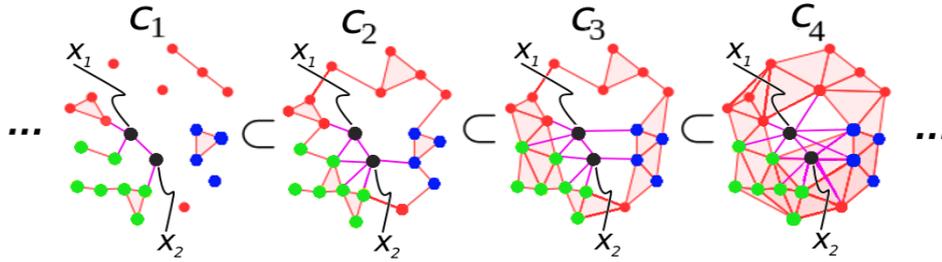
220 Now we define the algorithm.

221 3.2 Classification by using simplicial complexes and persistent 222 homology

223 The proposed method computes the predicted label list \hat{Y} corresponding to X according to
 224 definition 10 assuming T' is an incomplete association set. The entire process is summarized
 225 in Algorithm 1. We show a simplified example in Figure 4. The following section explains
 226 each step in detail.

228 **Algorithm 1 TDABC: TDA-Based Classification Algorithm**

229 **Require:** A training set $S \neq \emptyset$.
 230 A testing set $X \neq \emptyset$ to be classified.
 231 The incomplete association set T' .
 232 **Ensure:** A prediction list \hat{Y} of X by using T' .
 233 1: Obtain the filtration f' (see Figure 4) of the simplicial complex C constructed by using
 234 S and X (Algorithms 2 and 3)
 235 2: Obtaining the prediction list $\hat{Y} = (l_0, l_1, \dots, l_{|X|})$ where each $l_i \in \hat{Y}$ is the most reliable
 236 label corresponding to $x_i \in X$, with $0 \leq i < |X|$ by using f' filtration (Algorithm 4)
 237 3: **return** the prediction list \hat{Y} .



238 **Figure 4** TDABC (see algorithm 1) using a 4-level fragment of a filtration f' from a C simplicial
 239 complex. Two black points $x_1, x_2 \in X$ must be classified. The colored points are elements in S ,
 240 where each color represents different labels. By applying the Algorithm 1 in C_1, C_2, C_3, C_4 we get
 241 different labels for x_1 and x_2 depending on a selected C_i simplicial complex with $1 \leq i \leq 4$.

242 **3.2.1 Building the simplicial complex C and its filtration f' .**

243 We can build the simplicial complex C in two ways: using an incremental algorithm (Al-
 244 gorithm 2) or using a direct one (Algorithm 3). From Lemma 11 we know both algorithms
 245 are equivalents.

246 If we have a complex K from set S , we can use the incremental algorithm to insert in K
 247 every point $x \in X$. Each time we insert a point, we need to recompute or update the entire
 248 simplicial complex. In comparison to the direct algorithm, this strategy has the advantage
 249 of controlling resource consumption because it builds the simplicial complex step by step
 250 according to Equation 6. The direct strategy is useful when we need to classify a set and not
 251 just one point. In this case, it is better to build the complete simplicial complex once, in
 252 order to avoid updating the simplicial complex every time we insert a point.

284 **3.2.2 Obtaining the most reliable label**

285 Once we have a filtering f' from a simplicial complex C , we need to label all the elements
 286 of X , by using $\forall x \in X; \Upsilon(\{x\}, C)$ to assign the most voted label $l \in L$ during computation
 287 of $\Gamma(\{x\}, C)$ and $\Psi(\{x\}, C)$ according to definitions 6, ?? and 8.

288 As we know from Definition 13, C is the infinite level of the filtration f' . Theoretically the
 289 dimension of C can be $N - 1; N = |P|$ in case all the vertices are connected with all. We need
 290 to check $\forall x \in X, Lk(x)$ which contains the other elements of C . But, if $\exists x' \in X; x' \in Lk(x)$
 291 we need to analyze it recursively and so on until some $s \in S$ is found, which must contribute
 292 with its label. Resulting in a method with $O(|X| \cdot 2^N)$ is the worst case. Even with a table of
 293 dynamic programming to avoid recalculating the labels of the values twice, the computational

Algorithm 2 Incremental Construction of simplicial complex C

Require: A pre-existent filtration f from a non-empty simplicial complex K . A non-empty testing set X . An ξ increment value used to change the filtration level.

Ensure: f' a filtration of C simplicial complex (see Figure 4).

```

1: Let  $n \leftarrow |f|$  be the maximum level of filtration  $f$ 
2: Split  $X$  in  $n$  disjoint partitions  $\{X_0, X_1, \dots, X_n\}$ 
3:  $f' \leftarrow f$ 
4:  $i \leftarrow 0$ 
5: while  $i \leq n$  do
6:    $K_i \leftarrow \psi(f, i)$ , see definition 13
7:    $C_i \leftarrow K_i \uplus X_i$ 
8:   update  $f'$  with  $\psi(f', i) \leftarrow C_i$ 
9:    $i \leftarrow i + \xi$ 
10: end while
11: return  $f'$ 

```

Algorithm 3 Direct Construction of simplicial complex C

Require: A non-empty training set S . A non-empty testing set X . An ξ increment value used to change the filtration level. Let n be the maximum desired level of resulting filtration f' .

Ensure: f' a filtration of C simplicial complex (see Figure 4).

```

1:  $f' \leftarrow \{\emptyset\}$ 
2: Unify  $S$  and  $X$  by  $P \leftarrow S \cup X$ 
3:  $i \leftarrow 0$ 
4: while  $i \leq n$  do
5:   Build a simplicial complex  $C_i$  from  $P$  by using an  $i$  value [8, 1, 3, 2]
6:   Update  $i^{th}$  level of filtration  $f'$  with  $\psi(f', i) \leftarrow C_i$ 
7:    $i = i + \xi$ 
8: end while
9: return  $f'$ 

```

complexity would still be at least $O(2^N)$ that is required to initialize the table. If we could compute functions $\Psi(x, C)$ and $\Gamma(x, C)$ from definitions 6 y ?? without any change, it is very likely that the point $x \in X$ it would have all possible labels, due to the high degree of the analyzed simplices dimensions.

Two questions arise that address the calculation:

- (a) How do we know which simplices of C that contain some $x \in X$ are reliable and which are noise?
- (b) How can we make sure to assign a reliable label $\forall x \in X$?

We can use persistent homology to downsize the dimension of the problem. With persistence, we can get all topological invariants for all simplices dimension in f' . But we are just interested in working with q dimension; with $2 \leq q \ll N - 1$, which represents a complexity reduction of $O(|X| \cdot 2^{q+1})$ with q a constant value. Thus, we can say our method is linear in $|X|$. Those invariants are determined by persistence intervals with values (*birth*, *death*) (see definition 4).

For long life invariants (high *death* – *birth*) we are in the presence of a topological feature, but for short life, we have noise [13, 8, 29]. We need to find the level of filtration C_i that

326 **Algorithm 4** Labeling a testing set X set

327 **Require:** A filtration f' of a simplicial complex C .
328 A non-empty testing set X to classify.
329 The incomplete association set T' .

330 **Ensure:** A predicted labels list \hat{Y} of X according to definitions 10.

331 1: $D \leftarrow \text{ComputePersistentHomology}(f')$ [24, 9] where: $D = \{d_i \mid d_i = (\text{birth}, \text{death})\}$
332 2: Get a desired persistent interval $d \in D$ with $d \in \{\text{MaxInt}(D), \text{RandInt}(D), \text{AvgInt}(D)\}$
333
334 3: $i = d[\text{birth}]$
335 4: $C_i = \psi(f', i)$ see definition 13
336 5: $\hat{Y} \leftarrow \{\emptyset\}$
337 6: **while** $X \neq \emptyset$ **do**
338 7: $x \in X$
339 8: $l \leftarrow \Upsilon(\{x\}, C_i)$ see definitions 6, ??, 8
340 9: $\hat{Y} \leftarrow \hat{Y} \cup \{l\}$
341 10: $X \leftarrow X/\{x\}$
342 11: **end while**
343 12: **return** \hat{Y}

310 maximizes the number of topological features associated to each $x \in X$. At the same time,
311 we would like sub-complex C_i to have as many useful-simplices (see definition 5) as possible.
312 As a result, it is highly likely that we get a reliable label $\forall x \in X$.

313 As we have already described in Definition 4, a persistent interval $d = (\text{birth}, \text{death})$
314 represents the life time of a simplex set Q , in which every q-simplex was created and destroyed
315 on filtration levels $d[\text{birth}]$ and $d[\text{death}]$, respectively. Let D be the set of persistence interval
316 of q-simplices, and $d \in D$, then $\text{int}(d) = d[\text{death}] - d[\text{birth}]$. We define three ways to find
317 the desired persistent interval:

318 (a) The maximum persistent interval:

$$319 \quad d_m = \text{MaxInt}(D) = \text{Max}(\text{int}(d)); \forall d \in D. \quad (7)$$

320 (b) A persistent interval selected in a random way:

$$321 \quad d_r = \text{RandInt}(D) = \text{random}(D). \quad (8)$$

322 (c) The ceiling of persistent interval average:

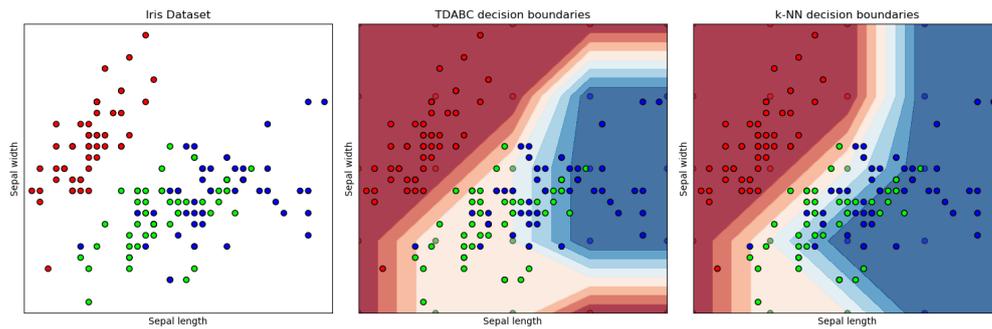
$$323 \quad d_a = \text{AvgInt}(D) = \lceil \text{Avg}(\text{int}(d)) \rceil; \forall d \in D. \quad (9)$$

324 Algorithm 4 executes the labeling process.

344 **4 Results**

345 We use the GUDHI library [19, 2, 1] to implement our TDA Based Classifier (TDABC), see
346 Algorithm 1. GUDHI is one of the most complete libraries for building simplicial complexes
347 and computing persistent homology [2, 21, 18, 1, 3].

348 To evaluate our proposed TDABC algorithm we use the Iris dataset [10]. The dataset
349 contains 3 classes of 50 instances each, where each class refers to a type of Iris plant. One
350 class is linearly separable from the other two; the latter are NOT linearly separable from each



357 ■ **Figure 5** Iris dataset (first two components)[left], decision boundaries from TDABC by using
 358 RandInt-function [center], and decision boundaries from k-NN [right].

351 other [10] (see Figure 5). L is the label set where $L = \{\text{“Setosa”}, \text{“Versicolor”}, \text{“Virginica”}\}$.
 352 Each sample in the Iris dataset is a 5-tuple, defined by:

$$353 \quad \text{5-tuple} = (\text{sepal_length}, \text{sepal_width}, \text{petal_length}, \text{petal_width}, \text{label})$$

354 We use the class of Iris plant as a predicted attribute. The left side of Figure 5 shows the
 355 distribution of the dataset (we choose the first two components for plotting purposes). We
 356 build a point set P using the first four components of every Iris dataset Sample.

359 4.1 Classifier Evaluation

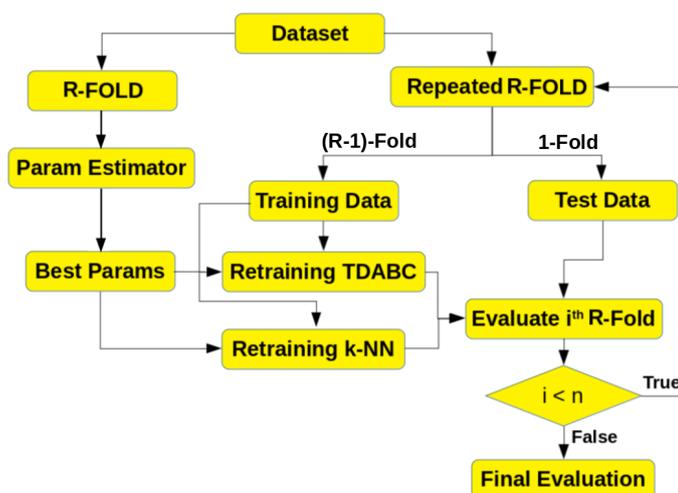
360 We create three different versions on our TDABC method depending on the applied selection
 361 function (see equations 7, 8 and 9): TDABC+MaxInt (TDABC-M), TDABC+RandInt
 362 (TDABC-R) and TDABC+AvgInt (TDABC-A). As a baseline algorithm, we choose the
 363 sci-kit learn k-Nearest Neighbors (k-NN) Classifier [22] to compare our results.

364 To make our evaluation more robust and avoid overfitting and underfitting we use Cross-
 365 Validation [22, 15, 20] method. The idea of Cross-Validation is to divide the data set P into
 366 equal pieces or folds. One piece of them is taken as the testing set X , and remaining folds
 367 as the set S . When we repeat this process, and folds of different sizes are generated, the
 368 method is called Repeated Cross-Validation.

369 Let R be the fold number in our Cross-Validation approach to avoid confusion with our
 370 use of k in k-NN. We decide to execute the repeated R-FOLD 5 times ($n=5$) and R will be
 371 one of (5, 10, 15, 20, 25) values, depending on the value of n , i.e.: $n=1$, $R=5$; $n=2$, $R=10$ and
 372 so on. For any value of R , we use $(R-1)$ -fold for a training set as S and the remaining fold as
 373 testing data as X in each iteration. The overall evaluation process is shown in Figure 6.

374 In machine learning algorithms it is common to use parameters whose values are set
 375 before the learning process begins. Those parameters are called hyper-parameters [22, 15,
 376 20]. By contrast, the values of other parameters are derived via training. In our evaluation,
 377 we need to find the best values for hyper-parameters k , and q (Param Estimator process in
 378 Figure 6). For k-NN we found $k=15$ as a good number of neighbors; we get it by using the
 379 hyper-parameter estimators from scikit-learn [22].

381 For the TDABC algorithm, we need to know the max simplex dimension q to control
 382 the VR-complex construction process. We fix $q = 3$ because the hardware limitation:16 GB
 383 RAM, Intel® Core™ i5 CPU 750 a 2.67GHz.



380 ■ **Figure 6** Repeated Cross Validation overall process to compare TDABC variants and k-NN.

384 4.2 Metrics for Classifiers Evaluation

385 In the interest of evaluating the performance of the proposed and baseline classifiers, we
 386 need to compute several metrics like: Accuracy (Acc), Precision (P), Recall (R), False
 387 Positive Rate (FPR), F1 measure ($F1$ or Armonic Measure of P and R), Mean Squared
 388 Error (MSE). As a graphical and general evaluation of a classifier we use the Confusion
 389 Matrix. With the exception of MSE , aforementioned metrics are defined by using the True
 390 Positives (TP), True Negatives (TN), False Positives (FP) and False Negatives (FN). Since
 391 the Iris dataset has multiple classes we consider TP_l, FP_l, TN_l, FN_l for each class $l \in L$. Let
 392 Y be the real label set $\forall x \in X$ (see definition 9). Let \hat{Y} be the predicted label set $\forall x \in X$
 393 computed by Algorithm 1, $n = |Y| = |\hat{Y}| = |X|$.

394 Therefore, we get each metric as the average of metric computation by each class $l \in L$,

395 according to equations from 10 to 19:

$$396 \quad TP_l = \sum_{i=1}^n \mathcal{I}(l = \hat{y}_i) \cdot \mathcal{I}(\hat{y}_i = y_i), \quad (10)$$

$$397 \quad FP_l = \sum_{i=1}^n \mathcal{I}(l = \hat{y}_i) \cdot \mathcal{I}(\hat{y}_i \neq y_i), \quad (11)$$

$$398 \quad TN_l = \sum_{i=1}^n \mathcal{I}(l \neq \hat{y}_i) \cdot \mathcal{I}(\hat{y}_i = y_i), \quad (12)$$

$$399 \quad FN_l = \sum_{i=1}^n \mathcal{I}(l \neq \hat{y}_i) \cdot \mathcal{I}(\hat{y}_i \neq y_i), \quad (13)$$

$$400 \quad Acc = \frac{1}{|L|} \cdot \sum_{l \in L} \frac{TP_l + TN_l}{TP_l + TN_l + FP_l + FN_l}, \quad (14)$$

$$401 \quad P = \frac{1}{|L|} \cdot \sum_{l \in L} \frac{TP_l}{TP_l + FP_l}, \quad (15)$$

$$402 \quad R = \frac{1}{|L|} \cdot \sum_{l \in L} \frac{TP_l}{TP_l + FN_l}, \quad (16)$$

$$403 \quad PR = \frac{1}{|L|} \cdot \sum_{l \in L} \frac{FP_l}{TN_l + FP_l}, \quad (17)$$

$$404 \quad F1 = 2 \cdot \frac{P \cdot R}{P + R}, \quad (18)$$

$$405 \quad MSE = \frac{1}{n} \cdot \sum_{i=1}^n (\hat{y}_i - y_i)^2. \quad (19)$$

406 We compute these 10 metrics [15, 4, 22] for any iteration of our repeated R-Fold strategy.

407 4.3 Comparison

408 Since $k \in (5, 10, 15, 20, 25)$ we execute each classifier a total number of $\sum_{i=1}^5 k_i = 75$ times.
409 We then compute the average of each metric and show the results in Table 1.

410	Name	Acc	P	R	FPR	MSE	F1
411	k-NN	0.97	0.92	0.89	0.04	0.10	0.91
412	TDABC-R	0.96	0.93	0.90	0.06	0.14	0.92
413	TDABC-A	0.95	0.91	0.88	0.05	0.29	0.90
414	TDABC-M	0.93	0.88	0.83	0.10	0.57	0.85

415 ■ **Table 1** Comparison of General Metrics by Classifier.

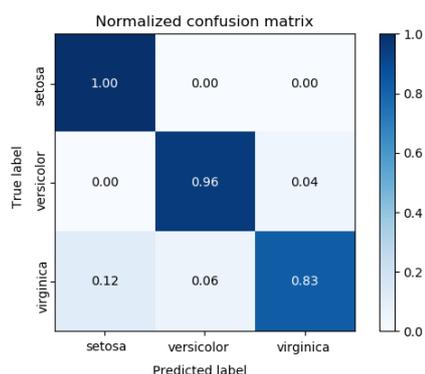
416 On the other hand, we compute the Overall Confusion Matrix (CM) for each classifier
417 concatenating all sets into one, resulting in two big sets of predicted labels and real labels:

$$418 \quad \hat{Y} = (\hat{Y}_1, \hat{Y}_2, \dots, \hat{Y}_{75}), \quad (20)$$

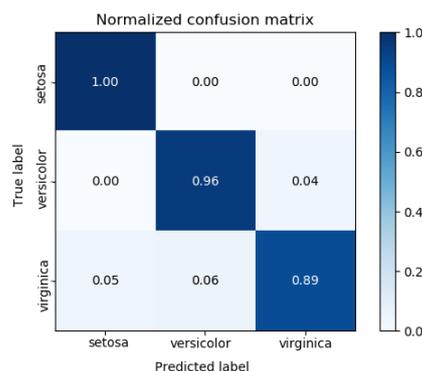
$$419 \quad Y = (Y_1, Y_2, \dots, Y_{75}). \quad (21)$$

420 where \hat{Y}_i is the predicted labels set and Y_i the real labels set, both resulting from i^{th} execution
421 of Repeated K-Fold. As we respect the relations between each predicted and real labels we
422 can easily calculate the CM for each classifier. The result is shown in the Figure 7.

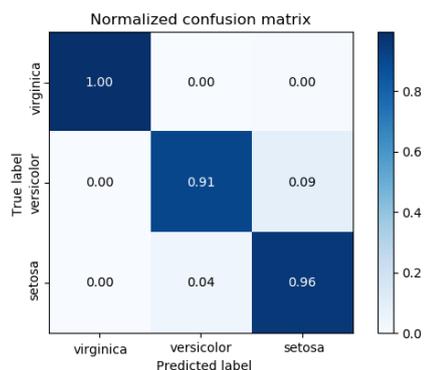
23:14 Classification based on Topological Data Analysis



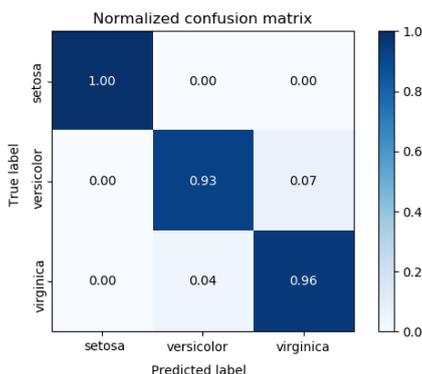
423 (a) TDABC-M.



423 (b) TDABC-A.



424 (c) TDABC-R.



424 (d) k-NN.

425 ■ **Figure 7** Confusion Matrices.

426 **5** Discussion

427 In section 3.2.1 we presented two algorithms to construct the simplicial complex using testing
 428 X and training S point sets. One is an incremental algorithm (Algorithm 2) and another is
 429 a direct one (Algorithm 3). Both are equivalents in their final result, as we have shown in
 430 Lemma 11. However, the incremental algorithm is recommended when a previous filtered
 431 simplicial complex was computed. Thus, a few values are classified if needed and you avoid
 432 recomputing all every time. This approach gives us the possibility of storing the filtered
 433 simplicial complex as a sort of knowledge representation (or topological database) for a
 434 specific dataset. The Direct algorithm is useful for training the TDABC with a new dataset
 435 or for classifying an entire point set at the same time.

436 The proposed voting system of assigning the most reliable label (see Definition 8), also
 437 offers the possibility of dealing with elements belonging to multiple classes at the same time.
 438 Perhaps in a more complex dataset, this possibility can be exploited in order to eliminate
 439 intra-class and extra-class variability.

440 The persistent homology is vital in Algorithm 4 to reduce the search space by taking
 441 advantage of topological features that are encoded inside selected persistent intervals.

442 Another interesting aspect is that the proposed method does not necessarily generate
 443 same-sized “neighborhoods” for each point to be classified as k-NN does. To get multi-sizes
 444 neighborhoods, no previous training is required, but k-NN needs a tuning process before
 445 inferring the best k hyper-param in order to get better results.

446 For the evaluation process, we use 10 metrics to cover as much as possible of analyzed
447 classification methods. In the end, we compute the overall confusion matrix for our three
448 classifiers (TDABC-R, TDABC-A, TDABC-M) and also for the k-NN classifier (see Figure 7).
449 Table 1 shows metric results in general terms, and our three TDABC have good classification
450 rates. However, the TDABC-R method is better than other TDBC analyzed methods, being
451 close to the k-NN classifier. On the other hand, TDABC-M is the worst, showing us that
452 more persistent topological features are not necessarily the best. Thus, more persistent
453 topological features could contain simplices with elements from several classes, resulting in
454 several errors during the labeling process with $\Psi(\sigma)$ function.

455 A more detailed analysis of the confusion matrices shows that for some classes the quality
456 of the proposed classification algorithms is equal to or surpasses k-NN. The TDABC-R
457 gets better results at classifying the Virginica class, see decision boundaries in Figure 5.
458 TDABC-A and TDABC-M are better at classifying Versicolor class and they are similar to
459 k-NN classifying Setosa with 100% accuracy. Therefore, in a detailed analysis of classification
460 per-class, our method is better than k-NN by detecting non-isolated classes.

461 **6** Conclusions

462 In this work, we develop a new classification method exclusively using Topological Data
463 Analysis tools. In particular, simplicial complexes and persistent homology. We also create
464 algorithms to build filtered simplicial complexes combining test and training sets. We also
465 create functions to label simplicial complexes and show how to use them in practice. The
466 use of persistent homology was determinant to reduce the complexity of the search space,
467 giving us a robust framework to understand data shapes and use it for classification. We
468 compare our proposed TDABC and its variants with k-NN obtaining similar rates. But, in
469 non-isolated classes, our methods were better. We think our method can be applied to solve
470 at least the same problems that k-NN solves. In this preliminary study, we experimentally
471 prove the possibility of classifying data by using TDA directly, with good accuracy rates.
472 Hopefully, this will open further research opportunities and of understanding data.

473 **7** Future Works

474 We have shown preliminary results about TDA-based classification. To our knowledge, this
475 is the first time TDA is directly used for classification. We still need to further explore new
476 data structures and algorithms to reveal and take advantage of the topological properties of
477 data. Furthermore, we want to process more complex datasets and evaluate the behavior of
478 our methods in the presence of noisy and high-dimensional data.

479 A new challenge is to find the lower simplicial complex dimension q needed to capture all
480 topological features from an arbitrary and high dimensional point set P . A lower q value is
481 useful for reducing the size of the searching space as much as possible and maintaining our
482 TDABC method in a linear or quasi-linear computation worst case. We do not know if there
483 is any relation between q and dataset sample dimensions. Finding this relation, we will be
484 able to understand how the optimum simplicial dimension q' can be found for any complex
485 dataset.

486 We want to explore the possibilities of our method in Feature Engineering as a data
487 imputation method in Missing Data Analysis (MDA). Finally, we want to better understand
488 the sort of problems and datasets in which our method could outperforms other methods.

8 References

- [1] Jean-Daniel Boissonnat and Karthik C. S. ‘An Efficient Representation for Filtrations of Simplicial Complexes’. In: *ACM Trans. Algorithms* 14.4 (2018), 44:1–44:21. URL: <https://doi.org/10.1145/3229146>.
- [2] Jean-Daniel Boissonnat, Karthik C. S. and Sébastien Tavenas. ‘Building Efficient and Compact Data Structures for Simplicial Complexes’. In: *Algorithmica* 79.2 (2017), pp. 530–567. DOI: 10.1007/s00453-016-0207-y. URL: <https://doi.org/10.1007/s00453-016-0207-y>.
- [3] Jean-Daniel Boissonnat and Clément Maria. ‘The Simplex Tree: An Efficient Data Structure for General Simplicial Complexes’. In: *Algorithmica* 70.3 (Nov. 2014), pp. 406–427. ISSN: 1432-0541. DOI: 10.1007/s00453-014-9887-3. URL: <https://doi.org/10.1007/s00453-014-9887-3>.
- [4] Michael W Browne. ‘Cross-Validation Methods’. In: *Journal of Mathematical Psychology* 44.1 (2000), pp. 108–132. ISSN: 0022-2496. DOI: <https://doi.org/10.1006/jmps.1999.1279>. URL: <http://www.sciencedirect.com/science/article/pii/S002224969912798>.
- [5] Gunnar Carlsson. ‘Topology and data’. In: *Bulletin of the American Mathematical Society* 46.2 (Jan. 2009), pp. 255–308.
- [6] Hamish Carr, Christoph Garth and Tino Weinkauf. *Topological Methods in Data Analysis and Visualization IV. Theory, Algorithms, and Applications*. Mathematics and Visualization. Gewerbestrasse, Switzerland: Springer International Publishing, 2017. ISBN: 9783319446844. DOI: 10.1007/978-3-319-44684-4.
- [7] Milka Doktorova and Afra Zomorodian. ‘Constructing Simplicial Complexes over Topological Spaces’. MsC. Thesis. MA thesis. New Hampshire, USA, 2012.
- [8] Herbert Edelsbrunner and John Harer. *Computational Topology - an Introduction*. Michigan, USA: American Mathematical Society, 2010. ISBN: 978-0-8218-4925-5. URL: <http://www.ams.org/bookstore-getitem/item=MBK-69>.
- [9] Herbert Edelsbrunner and Dmitriy Morozov. ‘Persistent homology: theory and practice’. In: *European Congress of Mathematics* (Jan. 2014), pp. 31–50. DOI: 10.4171/120-1/3.
- [10] Ronald A. Fisher. ‘The Use of Multiple Measurements in Taxonomic Problems’. In: *Annals of Eugenics* 7.7 (1936), pp. 179–188.
- [11] Kathryn Garside et al. ‘Topological data analysis of high resolution diabetic retinopathy images’. In: *PLOS ONE* 14.5 (May 2019), pp. 1–10. DOI: 10.1371/journal.pone.0217413. URL: <https://doi.org/10.1371/journal.pone.0217413>.
- [12] William Gasarch, Brittany Terese Fasy and Bei Wang. ‘Open Problems in Computational Topology’. In: *SIGACT News* 48.3 (Sept. 2017), pp. 32–36. ISSN: 0163-5700. DOI: 10.1145/3138860.3138867. URL: <http://doi.acm.org/10.1145/3138860.3138867>.
- [13] Robert Ghrist. ‘Barcodes: The persistent topology of data’. In: *BULLETIN (New Series) OF THE AMERICAN MATHEMATICAL SOCIETY* 45 (Feb. 2008), pp. 61–75. DOI: 10.1090/S0273-0979-07-01191-3.
- [14] Mostafa Hossny et al. ‘Driving behaviour analysis using topological features’. In: *2016 IEEE International Conference on Systems, Man, and Cybernetics (SMC)*. Budapest, Hungary: IEEE, Oct. 2016, pp. 003258–003263. DOI: 10.1109/SMC.2016.7844736.
- [15] Nathalie Japkowicz and Mohak Shah. *Evaluating Learning Algorithms: A Classification Perspective*. New York, NY, USA: Cambridge University Press, 2011. ISBN: 9780521196000.

- [16] Javier Lamar-Leon et al. ‘Persistent homology-based gait recognition robust to upper body variations’. In: *2016 23rd International Conference on Pattern Recognition (ICPR)*. Dec. 2016, pp. 1083–1088. DOI: 10.1109/ICPR.2016.7899780.
- [17] Ngoc-Khuyen Le et al. ‘Construction of the Generalized Czech Complex’. In: *2015 IEEE 81st Vehicular Technology Conference (VTC Spring)*. May 2015, pp. 1–5. DOI: 10.1109/VTCSpring.2015.7145759.
- [18] Clément Maria. ‘Algorithms and data structures in computational topology. (Algorithmes et structures de données en topologie algorithmique)’. PhD thesis. University of Nice Sophia Antipolis, France, 2014. URL: <https://tel.archives-ouvertes.fr/tel-01123744>.
- [19] Clément Maria et al. ‘The Gudhi Library: Simplicial Complexes and Persistent Homology’. In: *Mathematical Software – ICMS 2014*. Ed. by Hoon Hong and Chee Yap. Berlin, Heidelberg: Springer Berlin Heidelberg, 2014.
- [20] Tom M. Mitchell. *Machine learning, International Edition*. McGraw-Hill Series in Computer Science. McGraw-Hill, 1997. ISBN: 0071154671. URL: <http://www.worldcat.org/oclc/61321007>.
- [21] Nina Otter et al. ‘A roadmap for the computation of persistent homology’. In: *EPJ Data Science* 6.1 (Aug. 2017), p. 17. ISSN: 2193-1127. DOI: 10.1140/epjds/s13688-017-0109-5. URL: <https://doi.org/10.1140/epjds/s13688-017-0109-5>.
- [22] Fabian Pedregosa et al. ‘Scikit-learn: Machine Learning in Python’. In: *Journal of Machine Learning Research* 12 (Jan. 2012).
- [23] Lee M. Seversky, Shelvi Davis and Matthew Berger. ‘On Time-Series Topological Data Analysis: New Data and Opportunities’. In: *2016 IEEE Conference on Computer Vision and Pattern Recognition Workshops (CVPRW)*. June 2016, pp. 1014–1022. DOI: 10.1109/CVPRW.2016.131.
- [24] Vin de Silva, Dmitriy Morozov and Mikael Vejdemo-Johansson. ‘Persistent Cohomology and Circular Coordinates’. In: *Discrete & Computational Geometry* 45.4 (June 2011), pp. 737–759. ISSN: 1432-0444. DOI: 10.1007/s00454-011-9344-x. URL: <https://doi.org/10.1007/s00454-011-9344-x>.
- [25] Ann E. Sizemore et al. ‘The importance of the whole: Topological data analysis for the network neuroscientist’. In: *Network Neuroscience* 3.3 (2019), pp. 656–673. DOI: 10.1162/netn_a_00073. eprint: https://doi.org/10.1162/netn_a_00073. URL: https://doi.org/10.1162/netn_a_00073.
- [26] Primoz Škraba. ‘Persistent Homology and Machine Learning’. In: *Informatica. An International Journal of computing and Informatics* 42.2 (Jan. 2018), pp. 253–258.
- [27] Yuhei Umeda. ‘Time Series Classification via Topological Data Analysis’. In: *Transactions of the Japanese Society for Artificial Intelligence* 32 (May 2017), D-G72₁. DOI: 10.1527/tjsai.D-G72.
- [28] Vinay Venkataraman, Karthikeyan Natesan Ramamurthy and Pavan K. Turaga. ‘Persistent homology of attractors for action recognition’. In: *2016 IEEE International Conference on Image Processing, ICIP 2016, Phoenix, AZ, USA, September 25-28, 2016*. 2016, pp. 4150–4154. DOI: 10.1109/ICIP.2016.7533141. URL: <https://doi.org/10.1109/ICIP.2016.7533141>.
- [29] Afra Zomorodian. ‘Fast construction of the Vietoris-Rips complex’. In: *Computers & Graphics* 34 (3 June 2010), pp. 263–271.
- [30] Afra Zomorodian. *Topology for Computing*. Cambridge Monographs on Applied and Computational Mathematics. Cambridge University Press, 2009. ISBN: 0521136091.