# Persistent homology classification algorithm

**Mark Lexter D. De Lara** Corresp. 1, 2

1 Institute of Mathematical Sciences and Physics, College of Arts and Sciences, University of the Philippines Los Baños, College, Los Baños, Laguna, Philippines

2 Institute of Mathematics, University of the Philippines Diliman, Quezon City, Metro Manila, Philippines

Corresponding Author: Mark Lexter D. De Lara

Email address: mddelara@up.edu.ph

## ABSTRACT

Data classification is an important task in machine learning, used to solve problems in numerous settings. There are many classifiers, but none of the algorithms work best for all kinds of data, as implied by the no free lunch theorem. Topological data analysis is a rapidly growing field that deals with the shape of data. One primary tool in this field used to analyze the shape or topological properties of a dataset is persistent homology, a method based on algebraic topology for computing topological features of a space of points which persists across multiple resolutions. This study proposes a supervised learning and classification algorithm using persistent homology of training data classes in the form of persistence barcodes and diagrams to predict the output category of new observations. The developed algorithm was validated using real-world datasets and a synthetic dataset. The performance of the proposed classification algorithm on these datasets was compared to that of the most commonly used classifiers. Validation runs showed that the proposed persistent homology classification algorithm performed at par if not better than most of the classifiers considered.

## INTRODUCTION

Machine learning is a major branch of artificial intelligence. It deals with the study of computer systems and computer algorithms that can automatically learn and improve from experience without being explicitly programmed to do so. It focuses on the development of computer programs that can process data and give predictive analysis. Machine learning techniques are generally divided into three major categories, namely supervised learning, unsupervised learning, and reinforcement learning. In supervised learning, a system learns from a readily available training set of data with correctly labeled observations. One of the major tasks or problems addressed by supervised learning is classification.

Classification is the process of identifying, recognizing, grouping, and understanding new objects into categories/sub-populations (Alpaydin, 2014). A training dataset is composed of individual observations or n-dimensional data points which are split into an (n-1)-dimensional input vector often called features/explanatory variables, and into one-dimensional output vector/class/ label. These observations, also called instances, can be univariate, bivariate, or multivariate. These features, also called attributes, are quantifiable properties that can be categorical, ordinal, integer-valued, or real-valued. A classification algorithm, also called a classifier, is a procedure that implements classification tasks. Moreover, the term classifiers may also refer to the mathematical function that maps input features to an output category.

Classification algorithms have found many applications in the fields of computer vision, speech recognition, biometric identification, biological classification, pattern recognition, document classification, credit scoring, and many more. For instance, in medicine, the task of assigning a diagnosis to a given patient based on gathered features like age, gender, body mass index, presence of particular symptoms, etc., is a classification application. Classification problems can be categorized into binary classification or multi-class classification problems. Binary classification is the task of assigning an observation to exactly one of two categories, while multi-class classification is the process of assigning an instance to exactly one class out of more than two classes. Classification tasks tend to be harder in the presence of more than classes or more attributes.

**Deleted:** which

**Deleted:** has

**Deleted:** the ability to

**Deleted:** ,

**Deleted:** which

46   The study of classification algorithms is a vast field. Since the rise of artificial intelligence, numerous 47 classification algorithms have been developed. Several of these techniques can be used to solve binary classification 48   classification problems. Some algorithms are specially developed to solve binary classification problems, 49   while there are algorithms that can be used to solve binary and multi-class classification problems. Many 50 of these multi-class classifiers are extensions or modifications of one or more binary classifiers.

51   The no free lunch theorems proved by David Wolpert and William Macready in 1997 implies that 52   no learning or optimization algorithm that works best on all given problems (Wolpert and Macready, 53   1997). A classifier can be chosen depending on the type of data at hand. Since then, there had been so 54 many state-of-the-art classifiers that were developed. Some of the most commonly used classifiers are 55 logistic regression, multinomial logistic regression, Naive Bayes classifier, perceptron algorithm, linear 56 discriminant analysis, least squares support vector machines, quadratic classifiers, k-nearest neighbor 57 kernel density estimation, decision trees (random forests), and neural networks.

58 Many of these classifiers can be categorized as linear classifiers. A classification algorithm is a linear 59 classifier if it uses a linear function or linear predictor that assigns a score to each category $k$ based on 60 the dot product of a weight vector and the feature vector. The linear predictor is given by the score 61 functions, $Score(X_i, k) = \beta_k X_i$, where $X_i$ is the feature vector for the observation $i$, $\beta_k$ is the weight vector 62 corresponding category $k$. Observation $i$ is mapped by the linear predictor to the category $k$ with the 63 highest score function $\beta_k X_i$. Examples of linear classifiers include logistic regression, the perceptron 64 algorithm, support vector machines, and linear discriminant analysis (Yuan et al., 2012).

65   Data scientists employ techniques and theories drawn from many fields of mathematics, particularly 66   algebraic topology, statistics, information science, and computer science. In Mathematics, in particular, 67   there is a growing field called topological data analysis (TDA). It is an approach that uses tools and 68 techniques from topology to analyze datasets. In the past two decades, TDA has been applied in various 69 areas of science, engineering, medicine, astronomy, image processing, and biophysics.

70 One of the motivations in TDA is analyzing the shape of data and one of the main tools researchers use 71 is persistent homology (PH). PH is a method for computing topological features of a space of points which 72 persists across multiple resolutions (Carlsson, 2009),(Edelsbrunner and Harer, 2008),(Edelsbrunner and 73 Harer, 2010). It is based on the well-understood algebraic topology where invariant features can be derived 74 algebraically. These gathered invariant features are sensitive to small changes in the input parameters 75 which makes PH attractive to researchers who study qualitative features of data. PH involves representing 76 a point cloud by a filtered sequence of nested complexes, which are turned into novel representations 77 like barcodes and then interpreted statistically and qualitatively based on persistent topological features 78 which were gathered (Otter et al., 2017). A detailed discussion of pertinent information about the homology of 79 simplicial complexes and the process of computing persistent homology of a point cloud can be found in 80 the appendix.

81 Computation of PH has been applied in various areas including image analysis, shape comparison and 82 recognition, network analysis, computer visions, computational biology, oncology, chemical structures 83 and many more. Developments in the various aspects of computing PH have been increasing at a very rapid 84 rate. Various software were also developed to provide advanced and beginning practitioners platforms 85 to compute PH or develop new techniques in computing PH. These include JavaPlex, Perseus, Dipha, 86 Dionysus, jHoles, GUDHI, Rivet, Ripser, PHAT, R-TDA, and many more (Otter et al., 2017), (Pun et al., 87 2018).

88   This study is focused on the development of a supervised classification algorithm that mainly uses 89   persistent homologies of the datasets to solve classification problems. Persistent homology, which has been 90   around for only a decade has been getting so much attention in the past few years. Published works about 91   the fusion of these topics are quite new. Pun et al. (2018) published a survey of persistent-homology-based 92   machine learning algorithms and their applications. They presented a roadmap on how to use persistent 93 homologies to refine machine learning algorithms such as support vector machines, tree-based methods 94 and artificial neural networks. Their work was the inspiration in this study on how to extract topological 95   features based on persistence barcodes which resulted from

**Deleted:** which

**Deleted:** has

**Deleted:** homology

**Deleted:** its

**Deleted:** homology

**Deleted:**

computing data's persistent homology. In [96] their study, these features were considered as additional attributes to enhance machine learning algorithms. [97] While in this study, the topological features based on persistence barcodes/diagrams were directly used [98] and the main considerations in the proposed classifiers.

# PERSISTENT HOMOLOGY CLASSIFICATION ALGORITHM (PHCA)

The use of persistent homology in topological data analysis has been gaining attraction among researchers and data scientists. PH is mainly used to analyze the shape of a given dataset. A given point cloud undergoes a filtration process which turns it into a sequence of nested simplicial complexes. This is done by considering a finite number of increasing parameters and recording the sublevel sets that track changes in topological information. These changes can be documented in many ways, but the most popular ones are in terms of persistence barcodes or persistence diagrams. From these visualizations, the appearance and disappearances (birth and death) of intrinsic topological features like homology groups and Betti numbers are recorded and interpreted. The persisting duration (life span) of these topological features which are evident in the PH visualizations are essential in analyzing the qualitative and topological properties of data under study. PH has been used also to improve many machine learning algorithms. A list of these instances were mentioned and discussed by Pun et al. (2018). However, one of the main results of this study is the development of a supervised machine learning algorithm that mainly uses persistent homology of sets of data which can be used to solve classification problems.

Given a dataset or a point cloud composed of instances that belong to various classes, the first task is to divide the dataset into a training set and testing. Then, the goal is to analyze the dataset and develop a persistent-homology-based algorithm that will correctly identify the class to which each point in the testing set belongs to.

Consider a point cloud of size $M$ composed of $(n+1)$-dimensional data points. Suppose that in each point, the first $n$ entries are the attributes/features of the given point, and the $(n+1)$-th entry gives the class where the point belongs to. The $M$ points in the dataset are sorted into classes and each of the classes is split into a training set and testing set. For instance, in all the validation runs, we divide each of the classes into at least 80% training set and the remaining points into the testing set. Suppose there are $k$ classes and each class $i$, $i = 1,2,...,k$, is composed of $M_i$ points. Suppose also that in each class, there are $m_i$ points in the training set and $M_i - m_i$ points in the testing set. If $m$ is the sum of the $m_i$'s, then $m$ is the size of the training set, and $M - m$ is the size of the testing set.

Let $X$ be an $m \times (n+1)$ matrix in which rows represent the points in the training set. Similarly, let $Y$ be an $(M-m) \times (n+1)$ matrix containing the points in the testing set, the testing point cloud. Furthermore, let $X_i$ be an $m_i \times (n+1)$ matrix which contains the training set points belonging to class $i$. Call each of these matrices as training cloud for class $i$.

Before commencing the training, set the maximum dimension, denoted by *maxd*, that will be used in forming the Vietoris Rips complex filtration of the point clouds and computing the persistent homology of each of the training clouds. The parameter *maxd* is usually set to one or two during the validation runs. Validation runs show that these values of *maxd* are sufficient and the use of larger values of *maxd* will result to a longer computation time and may not be practical. Furthermore, there is also a need to set the maximum scales, denoted by *maxsc*. The scale here refers to the size of the epsilon balls to be considered in computing the persistent homology and topological features of the dataset. Preferably, the *maxsc* is set to be half the maximum distance between any two points in the point cloud.

After identifying the point cloud $X_i$ for class $i$, where $i$ goes from 1 to $k$ and setting *maxd* and *maxsc*, the algorithm may proceed to the following iterative steps.

Step 1. Training/Learning Stage For each $i$, $i \in \{1,...,k\}$, form the Vietoris Rips complex filtration for each point cloud $X_i$ for class $i$. Then, for each $i$, $i \in \{1,...,k\}$, compute the persistent homology of $X_i$, based on the Vietoris Rips complex filtration for each point cloud $X_i$. The result in computing the persistent homology of a point cloud is an $ntf \times 3$ matrix, where $ntf$ is the number of $d$-dimensional topological features that appear in the filtration. Denote this matrix by $P(X_i)$.

These topological features include the connected components, the loops, the voids, and so on. The

Deleted: s

Deleted:

Deleted: this

Deleted: which

Deleted: which

Deleted: a

Deleted: is

Deleted: the

Deleted: s

Deleted: Scale

Deleted: the

Deleted: s

145        number of topological features varies depending on the filtration. Let $tf_{i,j}$ be the $j$-th topological 146 feature of the point cloud $X_i$. The first column entries give the dimension of each of the topological

147 feature $tf_{i,j}$. Denote this by $dtf_{i,j}$, where $i = 1,...,k$ and $j = 1,...,ntf$. These entries take the values 148 0 for connected components, 1 for loops/holes, 2 for voids, and so on. The entries in the second 149 column give the birth time of each of the topological feature $tf_{i,j,}$ and the third column entries

150        gives the death time of each $tf_{i,j}$. Denote the birth time and death time of topological feature $tf_{i,j}$ 151 by $\beta_{i,j}$ and $\delta_{i,j}$ respectively. Visual presentation of each of the resulting persistent homology of a 152 given point cloud can be in the form of a persistence barcode or a persistence diagram.

153 **Step 2. Testing/Classification Stage.** For each of the $M-m$ data points in the testing set, identify the 154 class/category to which each data point belongs to.

155        Recall that $Y$ is an $(M-m)\times(n+1)$ matrix, where each row is a data point in the testing set. Let 156 $Y_j$ be the $j$-th row of $Y$ and the $j$-th data point in the testing set. Let the first $n$ entries of $Y_j$ be the 157 data point's attributes and the $(n+1)$-th entry be the data point's target class.

158 For each $j \in\{1,2,...,M-m\}$ and for each $i \in\{1,2,...,k\}$ append $Y_j$ to $X_i$ after the last row of 159 $X_i$. Name the resulting matrix $XY_{i,j}$. Perform filtration and PH computation on $XY_{i,j}$. That is,

160 compute $P(XY_{i,j})$. Record the change in topological features from $X_i$ to $XY_{i,j}$. Specifically, record 161 the change from $P(X_i)$ to $P(XY_{i,j})$. In this regard, consider two sets of point clouds, say point

162        cloud $A$ and point cloud $B$. Suppose there is an additional point $p$, to which we want to classify,

163        whether it belongs to point cloud $A$ or $B$. The proposed algorithm in this study will perform the 164 classification using topological features based on persistent homology. This technique is different 165 from the techniques used in the existing classifiers. Supposed that point $p$ is closer to point cloud $A$ 166 than point cloud $B$. Then, the persistent homology of $A\cup\{p\}$ possibly will have more topological 167 features compared to the persistent homology of $B\cup\{p\}$. Also, the birth of new topological features

168        will occur much earlier in $A\cup\{p\}$ and the death of some existing topological features may come earlier

169        in $A\cup\{p\}$.

170        With this phenomenon in mind, the terms in the score function, which measure the change in 171 topological features from $X_i$ to $XY_{i,j}$, are with reference to the following metrics.

172 (a) Let $\Omega_{i,j}$ be the difference of the sum of the entries of the first column of $P(X_i)$ from the sum 173 of the entries of the first column of $P(XY_{i,j})$.

174 (b) Let $\Phi_{i,j}$ be the difference of the sum of the entries of the third column of $P(X_i)$ from the 175 sum of the entries of the third column of $P(XY_{i,j})$.

176        (c) Let $\mu\Omega_{i,j}$ be the difference of the mean of the entries of the first column of $P(X_i)$ from the 177 mean of the entries of the first column of $P(XY_{i,j})$.

178        (d) Let $\mu\Phi_{i,j}$ be the difference of the mean of the entries of the third column of $P(X_i)$ from the 179 mean of the entries of the third column of $P(XY_{i,j})$.

180 (e) Let $AM_{i,j}$ be the sum over all $k$ of the absolute value of the difference of the mean of the 181 entries of the $k$-th column of $X_i$ and the mean of the entries of the $k$-th column of $XY_{i,j}$.

182        (f) Let $W_{i,j}$ be $p$-th Wasserstein distance of $P(X_i)$ from $P(XY_{i,j})$, where $p$ is set to 2.

The score function $Score(Y_j,i)$ is computed as

$$Score(Y_j,i) = -\Omega_{i,j} + \Phi_{i,j} - \mu\Omega_{i,j} + \mu\Phi_{i,j} + AM_{i,j} + W_{i,j}$$

**Deleted:** vary

**Deleted:** s

**Deleted:** s

**Deleted:** s

183    Finally, each data point $Y_j$ is assigned by the linear predictor to class *i* with the lowest score function 184 $Score(Y_j,i)$ over all *i*.

185    What follows is the pseudo-code for the persistent homology classification algorithm (PHCA).

**Deleted:** pseudo

# 186 EVALUATION METHODOLOGY

187 Classification is an instance of supervised learning. It is the task of identifying which of the categories 188 a new observation belongs to, based on a training set of data containing observations whose category

189    membership is known. Classifier is the term used to refer to the algorithm that implements the classification

190    and the mathematical function used by the classification algorithm to map an observation to a category. A

191    dataset is composed of (*n*+1)-dimensional data points, whose first *n* entries are called attributes of the 192 observation and the (*n*+1)-th entry is one of the *k* categories to which the observation belongs to. The

193    attributes can be real, integer, or categorical. The number of attributes, *n*, and the number of categories, *k*, 194 can be any fixed natural numbers. As the number of data points increases, or as the number of attributes 195 increases, the amount of computer time used to solve a classification problem also increases.

---

Algorithm 1 Persistent Homology Classification Algorithm

---

Require: $X_1, X_2, ...X_k, Y, maxd,$ and *maxsc*

Ensure: *Class*(Y) or *Class*($Y_j$) for each *j* procedure

TRAINING STAGE

∀*i* ∈{1,2,...,k}

  P($X_i$)←(*nt f*)×3 *matrix, a result of computing PH of $X_i$* end

procedure TESTING STAGE for *j* = 1 to *M* −*m* do for *i* = 1 to *k*

do

**Deleted:** procedure

        $XY_{ij}$ ← $X_i$ ∪{$Y_j$}
        P($XY_{ij}$)←(*nt f*)×3 *matrix, a result of computing PH of $XY_{ij}$*
        Compute for $\Omega_{i,j}, \Phi_{i,j}, \mu\Omega_{i,j}, \mu\Phi_{i,j}, AM_{i,j}, W_{i,j}$
        $Score(Y_j,i)=-\Omega_{i,j}+\Phi_{i,j}-\mu\Omega_{i,j}+\mu\Phi_{i,j}+AM_{i,j}+W_{i,j}$
        $Class(Y_j)$← arg min{$Score(Y_j,i)$}
                             ∀*i*

      end for

    end for

  end procedure

---

196 The classification algorithm developed in this study was validated by solving a number of classification 197 problems involving various classical validation datasets and a synthetic dataset. It should also be noted 198 that validation of the proposed algorithm in this study was implemented using *R* and the *R*-package TDA.

199    The different data used in the validation process were described in the following subsection.

**200 Validation Datasets.**

201    There were four datasets used in validating the proposed PHCA; three classical datasets and one synthetic 202 dataset. The number of classes per dataset is either two or three, while the number of attributes per dataset 203 ranges from two to seven.

204    1. Iris Plants Dataset

205    The Iris plant dataset created by Fisher (1936), available at the UCI Machine Learning Repository
206        (Dua and Graff, 2017), retrieved at https://archive.ics.uci.edu/ml/datasets/iris, one of the
commonly 207        used dataset in pattern recognition, is composed of 150 observations. The dataset
is divided into 208    3 categories or sub-populations, Iris Setosa, Iris Versicolour, and Iris Virginica.
Each category is 209 comprised of 50 data points. All of the 4 attributes of each data point, sepal
length, sepal width, 210        petal length, and petal width, are expressed in centimeters.

211    2. Wheat Seeds Dataset

212    The wheat seeds dataset was created by Charytanowicz et al. (2010) at the Institute of Agrophysics
       213        of the Polish Academy of Sciences in Lublin, available at the UCI Machine Learning
Repository
214 (Dua and Graff, 2017), and retrieved at https://archive.ics.uci.edu/ml/datasets/seeds. The dataset 215 is
composed of 210 observations which are divided equally into 3 categories: Kama, Rosa, and

**Deleted:** is

216        Canadian wheat variety. That is, there are 70 observations per category. Each data point is
217        characterized by seven attributes: area, perimeter, compactness, length of kernel, width of kernel,
218        asymmetry coefficient, and length of kernel groove. All of these parameters were real-valued and
       219        continuous.

220    3. Social Network Ads Dataset

221    The social network ads dataset was created by Raushan (2017) and retrieved at
       https://www.kaggle.
222    com/rakeshrau/social-network-ads/version/1. The dataset is composed of 400 data points. The 223
       observations were classified into two categories, whether a customer purchased a product (143) or
       224 not (257). Each data point has two attributes, age, and estimated salary. This classification task
is 225 considered as a bivariate classification problem.

226    4. Synthetic Dataset

227    The author created this dataset by generating 200 uniformly sampling points from each of the 228
       following figures, the circle defined by $x^2+y^2 = 25$, the sphere defined by $x^2+y^2+z^2 = 1$, and the

229        torus defined by $\left(3 - \sqrt{x^2 + y^2}\right)^2 + z^2 = 1$. The $x,y$, and $z$ coordinates of the 600 points served as
230        the attributes, and the category was assigned according to which figure the points belong to.

231    **Performance Measure.**
232        Measure of performance of the proposed PHCA were quantified and then compared with the
       performance 233 of major classification algorithms with respect to some validation datasets. The
       metrics used to evaluate 234 the methods were accuracy, sensitivity, and specificity. To compute for
       these metrics, the respective confusion 235 matrix for each method for the testing set was generated
       first. A confusion matrix is a table used to
236 describe the performance of a classification model on a set of test data for which the true values are 237
known. The confusion matrix gives the number of data points per class that are correctly predicted or 238
incorrectly predicted.

**Deleted:** which

239        For instance, consider a particular class, say $C_i$, among $k$ classes. Then, we can define the following

240        for each $i \in \{1,2,...,k\}$.

241        $TP_i$ is the number of true positives in class $C_i$, or the number of instances in $C_i$ which are predicted
       to

242        belong in $C_i$.

243        $TN_i$ is the number of true negatives in class $C_i$, or the number of instances outside $C_i$ which are
       predicted 244 to not belong in $C_i$.

245        $FP_i$ is the number of false positives in class $C_i$, or the number of instances outside $C_i$ which are
       predicted

246 to belong in $C_i$.

247 $FN_i$ is the number of false negatives in class $C_i$, or the number of instances in $C_i$ which are predicted to 248 not belong in $C_i$.

249 The three metrics per class $C_i$ are computed as follows

$$\text{Sensitivity of Class } C_i = \frac{TP_i}{TP_i + FN_i}$$

$$\text{Specificity of Class } C_i = \frac{TN_i}{TN_i + FP_i}$$

$$\text{Accuracy of Class } C_i = \frac{TP_i + TN_i}{TP_i + TN_i + FP_i + FN_i}$$

250 A high sensitivity prediction in Class $C_i$ implies that the reliability of predicting that an instance doesn't 251 belong to $C_i$ is high. However, predicting that an instance belongs to $C_i$ with high sensitivity is inconclusive.

252 On the other hand, the high specificity of prediction in Class $C_i$ implies that the reliability of predicting that an 253 instance belongs to $C_i$ is high. And, predicting that an instance doesn't belong to $C_i$ with high sensitivity is 254 inconclusive.

255 **Validation Procedure.**

256 The following procedure details the steps implemented to measure the performance of PHCA as compared 257 to other classification algorithms. These steps were performed for all of the four datasets.

258 1. Consider the dataset as a point cloud $X$. Divide it into 2 parts, training set, and testing set. For all 259 the validation runs, we have split the dataset to at least 80% training set and the remainder to testing 260 set.

261 2. Solve the classification problem using the proposed PHCA and each of the five algorithms: Linear 262 discriminant analysis (LDA), Classification and Regression Trees (CART), K-Nearest Neighbors 263 (KNN), Support Vector Machine (SVM), and Random Forest (RF). Depending on the algorithm 264 used, utilize the training set and classify each point in the testing set. Information about the nature of

265 these classifiers, including examples and program codes, are available in Subasi (2020), Stanimirova 266 et al. (2013), Breiman et al. (1984), Loh (2011), Neath and Johnson (2010), Cortes and Vapnik 267 (1995), Ho (1995), and Ho (1998).

268 3. Construct the confusion matrix per classification algorithm.

269 4. Compute the performance of each classification algorithm in terms of accuracy, and sensitivity, and 270 specificity per class.

## RESULTS AND DISCUSSION

272 Presented here are the performance of the proposed PHCA and the five major classification algorithms
273 in solving four classification problems. Program codes written in R which implements PHCA, LDA,
274 CART, KNN, SVM, and RF can be found on https://github.com/mlddelara/PHCA. There is a section for 275 the discussion of validation results for each of the classification tasks. Presented in each section are the
276 persistence diagrams and the persistence barcodes of the training sets. Recall that PHCA works in a way
277 that a data point in the testing set will be classified under a class if its inclusion in the particular class' 278 training set results to the least change in the persistence diagram or persistence barcode of the training set 279 with the additional data point.

280 **Iris Plants Dataset.**

281 The dataset is comprised of 50 data point from each of the three types of iris plant, namely, Iris Setosa,
282 Iris Versicolour, and Iris Virginica. Each data point is composed of four features and a class label. For
283   each class, ten data points were set aside to be part of the testing set and the remaining forty points
were 284        collected as the training set per class.
285      Figures 1, 2, and 3 shows the persistence diagram and persistence barcode of the respective training
286      sets. These are the representations of computing the persistent homology of each of the training set
per 287 class.

**Figure 1.** Peristence Diagram and Barcode for the Iris Setosa (Class 1) Training Set

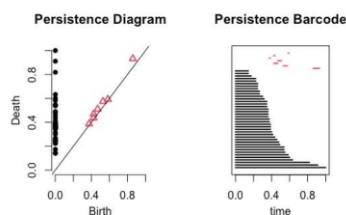**Figure 2.** Peristence Diagram and Barcode for Iris Versicolour (Class 2) Training Set

**Figure 3.** Peristence Diagram and Barcode for Iris Virginica (Class 3) Training Set

288 Table 1 shows the performance of PHCA and the five major classification algorithms in terms of 289
accuracy, sensitivity per class, and specificity per class. PHCA ranked third in terms of accuracy. That
290 is, of 30 testing data points, only one was wrongly classified. SVM performed equivalently with PHCA, 291
while CART and RF performed poorer with 2 mistakes each. On the other hand, LDA and KNN performed
292 perfectly for this problem.

| Classifier | Accuracy | Sensitivity per class | Specificity per class |
|---|---|---|---|

|        | Class 1 | Class 2 | Class 3 | Class 1 | Class 2 | Class 3 |
|--------|---------|---------|---------|---------|---------|---------|
| LDA    | 100%    | 100%    | 100%    | 100%    | 100%    | 100%    | 100% |
| CART   | 93.33%  | 100%    | 100%    | 80%     | 100%    | 90%     | 100% |
| KNN    | 100%    | 100%    | 100%    | 100%    | 100%    | 100%    | 100% |
| SVM    | 96.67%  | 100%    | 100%    | 90%     | 100%    | 95%     | 100% |
| RF     | 93.33%  | 100%    | 100%    | 80%     | 100%    | 90%     | 100% |
| PHCA   | 96.67%  | 100%    | 90.91%  | 100%    | 100%    | 100%    | 95.24% |

| Number of Data Points: | 150 | Number of Classes: | 3 |
|---|---|---|---|
| Training Set Size: | 120 | Number of Attributes: | 4 |
| Testing Set Size: | 30 | | |

**Table 1.** Result of classifying the Iris dataset using the six classifiers

**Wheat Seeds Dataset.**
The dataset is comprised of 70 data points for each of the three types of wheat varieties, namely, Kama, Rosa, and Canadian. Each of the data points has seven attributes and a class label. For each class, there are 14 testing data points and 56 training data points.

The persistence diagram and persistence barcode of the respective training set per class was computed and represented in Fig. 4, Fig. 5, and Fig. 6.
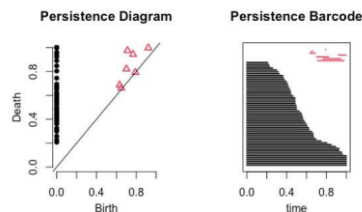


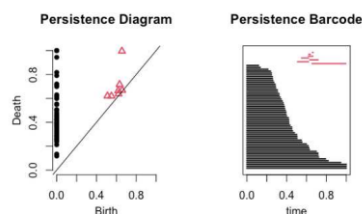**Figure 4.** Peristence Diagram and Barcode for Kama Variety (Class 1) Training Set



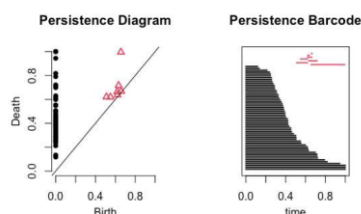**Figure 5.** Peristence Diagram and Barcode for Rosa Variety (Class 2) Training Set

**Figure 6.** Peristence Diagram and Barcode for Canadian Variety (Class 3) Training Set

₂₉₉ Table 2 shows the performance of PHCA and the five major classification algorithms in terms of ₃₀₀ accuracy, sensitivity per class, and specificity per class. PHCA got the highest accuracy, together with ₃₀₁ RF and SVM. These algorithms wrongly classified only one data point among 42 testing data points. ₃₀₂ Moreover, PHCA got the highest sensitivity and specificity for each of the classes. On the other hand, ₃₀₃ CART performed the worst in terms of accuracy which wrongly classified three data points.

| Classifier | Accuracy | Sensitivity per class | | | Specificity per class | | |
|---|---|---|---|---|---|---|---|
| | | Class 1 | Class 2 | Class 3 | Class 1 | Class 2 | Class 3 |
| LDA | 95.24% | 92.86% | 100% | 92.86% | 96.43% | 100% | 96.43% |
| CART | 92.86% | 92.86% | 100% | 85.71% | 92.86% | 100% | 96.43% |
| KNN | 95.24% | 92.86% | 100% | 92.86% | 96.43% | 100% | 96.43% |
| SVM | 97.62% | 100% | 100% | 92.86% | 96.43% | 100% | 100% |
| RF | 97.62% | 100% | 100% | 92.86% | 96.43% | 100% | 100% |
| PHCA | 97.62% | 100% | 100% | 93.33% | 96.55% | 100% | 100% |

| Number of Data Points: | 210 | Number of Classes: | 3 |
|---|---|---|---|
| Training Set Size: | 168 | Number of Attributes: | 7 |
| Testing Set Size: | 42 | | |

**Table 2.** Result of classifying the Wheat Seeds dataset using the six classifiers

₃₀₄ **Social Network Ads Dataset.**

₃₀₅ The dataset is comprised of uneven number of observations per class. There are 143 data points for class ₃₀₆ 1 and 257 data points for class 2. Each of the data points has two attributes and a class label. The former ₃₀₇ represents observations from customers who purchased a product. There are a total of 80 testing data ₃₀₈ points and 320 training data points.

₃₀₉ The persistence diagram and persistence barcode of the respective training set per class was computed ₃₁₀ and shown in the Fig. 7 and Fig. 8.
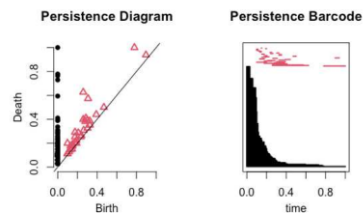
**Figure 7.** Peristence Diagram and Barcode for Purchaser (Class 1) Training Set
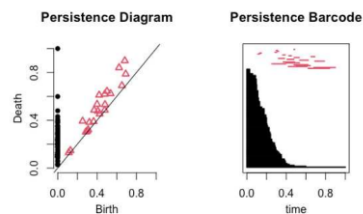


**Figure 8.** Peristence Diagram and Barcode for Non-purchaser (Class 2) Training Set

311 Table 3 shows the performance of PHCA and the five major classification algorithms in terms of 312 accuracy, sensitivity, and specificity. PHCA ranked third in terms of accuracy. It got 100% sensitivity, but 313 lower specificity at 90.91%. SVM performed equivalently with PHCA, in terms of accuracy. LDA and 314 KNN got 100% accuracy, but RF and CART got the lowest accuracy of 93.33%.

| Classifier | Accuracy | Sensitivity | Specificity |
|---|---|---|---|
| LDA | 86.08% | 90.20% | 78.57% |
| CART | 87.34% | 86.27% | 89.29% |
| KNN | 82.28% | 92.16% | 64.29% |
| SVM | 86.08% | 88.24% | 82.14% |
| RF | 86.08% | 90.20% | 78.57% |
| PHCA | 82.72% | 91.30% | 71.43% |
| Number of Data Points: | 400 | | |
| Training Set Size: | 181 | | |
| Testing Set Size: | 119 | | |
| Number of Classes: | 2 | | |
| Number of Attributes: | 2 | | |

**Table 3.** Result of classifying the Social Network Ads dataset using the six classifiers

315 **Synthetic Dataset.**

The dataset is comprised of 600 data points. There are three classes with 200 data points per class. Each 317 of the data points has three attributes and a class label. For each of the three classes, there are 40 testing 318 data points and 160 training data points.

319      The persistence diagram and persistence barcode of the respective training set per class was computed 320 and represented in Fig. 9, Fig. 10, and Fig. 11.
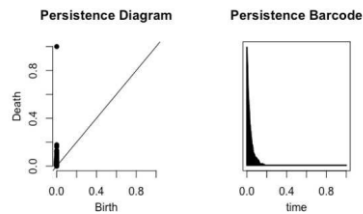


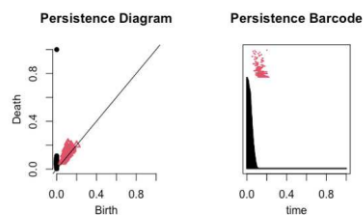**Figure 9.** Peristence Diagram and Barcode for Circle (Class 1) Training Set



**Figure 10.** Peristence Diagram and Barcode for Sphere (Class 2) Training Set
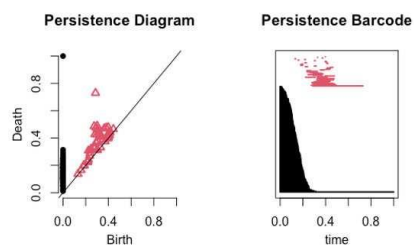


**Figure 11.** Peristence Diagram and Barcode for Torus (Class 3) Training Set 321 Table 4 shows the performance of PHCA and the five major classification algorithms in terms of 322

accuracy, sensitivity, and specificity. PHCA, together with CART, KNN, SVM, and RF, performed

323 perfectly with 100% accuracy, sensitivity per class, and specificity per class. While LDA got a low 324 accuracy of 93.33%.

| Classifier | Accuracy | Sensitivity per class | | | Specificity per class | | |
|---|---|---|---|---|---|---|---|
| | | Class 1 | Class 2 | Class 3 | Class 1 | Class 2 | Class 3 |
| LDA | 93.33% | 100% | 82.50% | 97.50% | 91.25% | 98.75% | 100% |
| CART | 100% | 100% | 100% | 100% | 100% | 100% | 100% |
| KNN | 100% | 100% | 100% | 100% | 100% | 100% | 100% |
| SVM | 100% | 100% | 100% | 100% | 100% | 100% | 100% |
| RF | 100% | 100% | 100% | 100% | 100% | 100% | 100% |
| PHCA | 100% | 100% | 100% | 100% | 100% | 100% | 100% |

| Number of Data Points: | 600 | Number of Classes: | 3 |
|---|---|---|---|
| Training Set Size: | 480 | Number of Attributes: | 3 |
| Number of Data Points: | 120 | | |

**Table 4.** Result of classifying the Synthetic dataset using the six classifiers

325 Validation of the performance of PHCA was done by comparing its performance in solving four
326 classification problems against the respective performances of the five major classification algorithms 327 in solving the same problems. The four validation datasets are comprised of a varying number of data 328 points, number of classes and number of attributes per observation. In terms of accuracy, sensitivity, and 329 specificity, PHCA and all the benchmark algorithms, excluding LDA, ranked first in two of four validation 330 data sets. However, only PHCA and SVM faired well in all four classification problems. All the other 331 algorithms had the worst accuracy, sensitivity, and specificity in at least one of the problems. CART has 332 the worst performance in solving the Iris dataset and Seeds dataset. LDA has the worst performance 333 in solving the synthetics dataset. And, KNN and RF have the worst performance in solving the Social 334 Network Ads dataset and Iris dataset, respectively.

335 These validation results do not imply that PHCA is better than any of the other major classification 336 algorithms. But, these results are just evidences to the no free lunch theorem which implies that no
337 learning algorithm works best on all given problems. Moreover, these validation runs imply that PHCA 338 can be at par or even better than some other classifiers in solving some particular classification problems.

339 What sets PHCA apart from the well-known machine learning classifiers is that it is non-parametric,
340 but at the same time a linear classifier. It is a non-parametric algorithm in the sense that it does not restrict 341 the data to follow a particular distribution nor fix the number of datasets' parameters for the algorithm
342 to work. PHCA works by assigning topological attributes from persistent homology of training data 343 points per classes and uses this as the parameters needed for a linear classifier which the algorithm uses to 344 classify new points. The referred topological attributes include the dimension, birth time, and death time 345 of topological features of the different training datasets and classes, and the Wasserstein distance between 346 classes.

347 **CONCLUSIONS**

348 The main result of this study was the development of PHCA, a non-parametric but linear classifier which 349 utilizes persistent homology, a major and very powerful TDA tool. Classification tasks are major concerns 350 in the field of machine learning which is why solving these kinds of problems has been a widely studied 351 discipline. The proliferation of the various classification algorithms is further fueled by

**Deleted:** is

the fact implied ₃₅₂     by the no free lunch theorem which implies that there is no single best algorithm which can be used to ₃₅₃    solve all types of classification problems.

₃₅₄PHCA was validated in this study by using it to solve four different classification problems with ₃₅₅ varying sizes, number of classes, and number of attributes. PHCA's performance per problem-based ₃₅₆ on accuracy, sensitivity, and specificity was measured and compared with the performance of five other ₃₅₇ well-known classifiers. The validation runs show that PHCA can perform well, or even better, than some ₃₅₈     of the major supervised machine learning classifiers, in solving particular classification tasks. Moreover, ₃₅₉    this validation activity does not imply that PHCA works better than other machine learning algorithms, ₃₆₀    but this exposition shows that PHCA can work in solving some classification problems.

₃₆₁Validation in this study was limited to relatively small problems which are restricted by the computers ₃₆₂ used in this study. PHCA can be further validated by considering larger problems and by using more ₃₆₃ powerful computers which can solve problems with higher dimensions. These future researches could ₃₆₄ test whether PHCA can still perform at par with or better than other classifiers. Furthermore, various ₃₆₅ improvements may be imposed on the proposed classification algorithm in this study by considering other ₃₆₆ topological attributes or by considering persistent homology representations other than barcodes and ₃₆₇ diagrams. Recent improvements and modifications on the computation of persistent homology may also ₃₆₈ be adapted to possibly improve the performance of PHCA. PH computations and the validation of the ₃₆₉ proposed algorithm were implemented using $R$ and TDA package in $R$. It should be noted that there are ₃₇₀    other platforms and solvers which can be used, like JavaPlex, Perseus, Dipha, Dionysus, jHoles, GUDHI, ₃₇₁    Rivet, Ripser and PHAT, which offer some variations in the the way PH can be computed. Indeed, this ₃₇₂     study has opened a lot of research opportunities which can be explored by mathematicians, data scientists, ₃₇₃ topologists, and computer programmers.

₃₇₄**ACKNOWLEDGMENTS**

₃₈₀**REFERENCES**

₃₈₁Alpaydin, E. (2014). *Introduction to Machine Learning, third edition*. Adaptive Computation and Machine ₃₈₂Learning series. MIT Press.

₃₈₃    Breiman, L., Friedman, J. H., Olshen, R. A., and Stone, C. J. (1984). *Classification and regression trees*. ₃₈₄    Wadsworth and Brooks/Cole Advanced Books and Software.

₃₈₅    Carlsson, G. (2009). Topology and data. *Bull Am Math Soc*, 46.

₃₈₆    Charytanowicz, M., Niewczas, J., Kulczycki, P., Kowalski, P., Lukasik, S., and Zak, S. (2010). A complete ₃₈₇ gradient clustering algorithm for features analysis of x-ray images.

₃₈₈    Cortes, C. and Vapnik, V. N. (1995). Support-vector networks. *Machine Learning*, 20(3):273–297.

₃₈₉    Dua, D. and Graff, C. (2017). UCI machine learning repository.

₃₉₀    Edelsbrunner, H. and Harer, J. (2008). Persistent homology — a survey.

₃₉₁    Edelsbrunner, H. and Harer, J. (2010). *Computational topology: an introduction*. Am. Math. Soc., ₃₉₂    Providence.

₃₉₃    Fisher, R. (1936). The use of multiple measurements in taxonomic problems. *Annual Eugenics*, 7(Part ₃₉₄ II):179–188.

₃₉₅    Ho, T. K. (14–16 August 1995). Random decision forests. *Proceedings of the 3rd International Conference ₃₉₆ on Document Analysis and Recognition, Montreal, QC,*, page 278–282.

[397]    Ho, T. K. (1998). The random subspace method for constructing decision forests. *IEEE Transactions on* [398]   *Pattern Analysis and Machine Intelligence*, 20(8):832–844.

[399]    Loh, W. (2011). Classification and regression trees. *Wiley Interdisciplinary Reviews: Data Mining and* [400]   *Knowledge Discovery*, 1(1):14–23.

[401] Neath, R. and Johnson, M. (2010). Discrimination and classification. In Peterson, P., Baker, E., and [402] McGaw, B., editors, *International Encyclopedia of Education (Third Edition)*, pages 135–141. Elsevier, [403] Oxford, third edition edition.

[404] Otter, N., Porter, M. A., Tillmann, U., Grindrod, P., and Harrington, H. A. (2017). A roadmap for the [405] computation of persistent homology. *EPJ Data Science*, 6(1):17.

[406]    Pun, C. S., Xia, K., and Lee, S. X. (2018). Persistent-homology-based machine learning and its applica-

[407]    tions – a survey.

[408]    Raushan, R. (2017). Social network ads, version 1.

[409]    Stanimirova, I., Daszykowski, M., and Walczak, B. (2013). Robust methods in analysis of multivariate [410]food chemistry data. *Chemometrics in Food Chemistry*, page 315–340.

[411]    Subasi, A. (2020). Machine learning techniques. In *Practical Machine Learning for Data Analysis Using* [412] *Python*, pages 92–202. Academic Press.

[413] Wolpert, D. H. and Macready, W. G. (1997). No free lunch theorems for optimization. *IEEE Transactions* [414] *on Evolutionary Computation*, 1(1):67–82.

[415]    Yuan, G., Ho, C., and Lin, C. (2012). Recent advances of large-scale linear classification. *Proceedings of* [416] *the IEEE*, 100(9):2584–2603.