

# A multi prototype classification algorithm and its application to multi class diagnostics

Mario Ziller

This paper introduces a novel, universal distance-based classification procedure. It is based on a simple geometric model. Considering all objects as points in a metric space, a class is imagined as covered by potentially differentsized hyperspheres, the centres of which are referred to as prototypes. The radii of the hyperspheres are individually optimised by a generalised ROC-analysis. For the approximate solution of the entire discrete optimisation problem, a greedy algorithm was developed and implemented in R. It runs in  $O(k^2 \cdot n^2 \cdot \log(n))$  time where  $k$  is the number of prototypes to be selected and  $n$  the number of training objects. For application to multi class problems, one against all approach is performed. The diagnostic decision is finalised for that class of maximum positive predictive value when in doubt. Objects not recognised as a member of any of the classes are assigned to an additional residual class. The performance of the classification system presented is demonstrated on various data examples, and in comparison with other methods.

# 1 A multi prototype classification algorithm 2 and its application to multi class 3 diagnostics

4 **Mario Ziller<sup>1</sup>**

5 <sup>1</sup>**Biomathematics Working Group, Friedrich-Loeffler-Institut, Federal Research Institute  
6 for Animal Health, Greifswald - Insel Riems, Germany**

## 7 **ABSTRACT**

This paper introduces a novel, universal distance-based classification procedure. It is based on a simple geometric model. Considering all objects as points in a metric space, a class is imagined as an overlap of potentially different-sized hyperspheres, the centres of which are referred to as prototypes. The radii of the hyperspheres are individually optimised by a generalised ROC-analysis. For the approximate solution of the entire discrete optimisation problem, a greedy algorithm was developed and implemented in R. It runs in  $\mathcal{O}(k^2 \cdot n^2 \cdot \log(n))$  time where  $k$  is the number of prototypes to be selected and  $n$  the number of training objects.

For application to multi class problems, one against all approach is performed. The diagnostic decision is finalised for that class of maximum positive predictive value when in doubt. Objects not recognised as a member of any of the classes are assigned to an additional residual class. The performance of the classification system presented is demonstrated on various data examples, and in comparison with other methods.

Keywords: Prototype classifier, Global distance-based classification, ROC analysis, Greedy algorithm,  
9 Multi class diagnostics

## 10 **1 INTRODUCTION**

11 Many diagnostic decisions in medicine, biology, and far beyond are substantiated on large pattern data  
12 bases e.g. of DNA-sequences, mass spectra, or others. These data bases often include many predefined  
13 diagnostic classes (Wong et al., 2010; De Bruyne et al., 2011; Karger et al., 2011; Hong et al., 2008). The  
14 classification is usually done by comparison based on distances or similarity scores. Screening for similar  
15 proof-samples by local distance-based methods makes use of the nearest neighbour principle. This works  
16 well but may be computationally intensive (Tan et al., 2005). For classification assignment, the distances

17 of an unknown object to all of the data base objects must be assessed. Therefore, it may become ineligible  
18 for application in high-throughput diagnostics e.g. of emerging diseases in case of a pandemic situation.

19 For comparable problems, a global approach may be more useful. Its application to large data sets  
20 operates much faster. There are various global distance-based methods like centroid- or medoid-based  
21 classification, etc. (Breiman et al., 1984; Tan et al., 2005; Gordon, 1999; Hastie et al., 2009). They use  
22 one or more class representatives for comparison. One can imagine that in general class members should  
23 have lower distances to other class members than to non-members. Therefore, the choice of the distance  
24 measure is of great importance for the correctness of subsequent class-assignments. Simple k-means or  
25 k-medoids models consider only one representative of each class. They may be effective if all classes  
26 are spherically shaped (Tan et al., 2005; Hastie et al., 2009). Plain assignments to the nearest centroids  
27 or medoids, which other classifiers like e.g. neural gases (Martinetz and Schulten, 1991) make use of,  
28 however, would require either classes of nearly equally sized hyperspheres or the selection of possibly a  
29 great many representatives.

30 All this might be improper in real practical applications. In this context, a more general approach  
31 (Ziller et al., 2011) was developed. In this, a class is considered being covered by several different-sized  
32 hyperspheres. Hence, more than one centre of each class and different radii are inspected. For the  
33 optimum sizing of those hyperspheres receiver operating characteristic (ROC)-analysis is utilised, which  
34 is a well-established standard framework for evaluating diagnostic tests and assessing suitable cutoffs  
35 (Hanley and McNeil, 1982; Greiner et al., 2000). For the classification task under consideration, this  
36 fundamental idea was generalised for application to multi-core models.

37 In this paper, a novel multi prototype classifier is outlined. The basic classification algorithm for single  
38 class problems and its generalisation to multi class problems are expounded. Performance reliability of  
39 the classifier is demonstrated on various data sets in comparison with other well-established classification  
40 tools.

## 41 **2 MULTI PROTOTYPE CLASSIFICATION**

### 42 **2.1 Theoretical framework**

#### 43 ***Mathematical model***

44 For a short mathematical outline, let all objects be considered as points in a given metric space. Any class  
45 to be investigated is modelled as a union of various, potentially different-sized hyperspheres which may  
46 overlap. The centres of the hyperspheres represent reference objects of the class, henceforth referred to  
47 as 'prototypes'. The radii of the hyperspheres are used as thresholds, hereafter referred to as 'cutoffs',  
48 for the diagnostic decision. They are individually optimised by a generalised ROC-analysis which all

49 other hyperspheres were fixed in. An object is assigned to the class under consideration if its distance  
 50 to at least one of the prototypes is not greater than the respective cutoff. In this approach only distances  
 51 between the objects are given. Feature vectors or other information are not available. No additional model  
 52 assumptions, e.g. about the distribution of point distances or about cluster morphologies, i.e. a radially  
 53 decreasing density of objects within the class, were made.

Let  $O$  be the set of  $n$  objects considered,  $C \subset O$  the subset of  $m$  class members, and  $d$  a metric distance function on  $O$ . The class should be characterised by the prototypes  $x_i \in C, i = 1, \dots, k$  and the related cutoffs  $f_i \in R \geq 0$ . The number  $k$  of prototypes is arbitrary but fixed and must have been selected before the analysis starts. The membership of an object  $y \in O$  to the class  $C$  was then modelled by the decision rule

$$y \in \hat{C} \Leftrightarrow \exists i \in 1, \dots, k : d(y, x_i) \leq f_i.$$

The set of all predicted class members is therefore

$$\hat{C} = \bigcup_{i=1}^k \hat{C}_i = \bigcup_{i=1}^k \{y \in O : d(x_i, y) \leq f_i\}.$$

#### 54 **ROC-Analysis**

55 ROC-analysis is a well-established and for this purpose generalised method of estimating thresholds  
 56 for effective dichotomisation. This procedure represents an essential part of the investigated algorithm  
 57 and will be massively applied in it. Cutoff thresholds for diagnostic decisions can be determined by  
 58 classical ROC-analyses (Hanley and McNeil, 1982; Greiner et al., 2000). For a given feature function  
 59  $f$ , a threshold  $f^*$  should be determined which makes  $f(x) \leq f^*$  a good classifier, i.e. according to the  
 60 declared notation, an object  $x \in O$  will be assigned to class  $C$  if and only if  $f(x) \leq f^*$ . In the described  
 61 classification algorithm, the distance to a fixed prototype is used as feature-function. The best cutoff  $f^*$   
 62 can be found by optimising a criterion function  $g$ .

As usual, sensitivity  $Se$  is defined as the proportion of correctly predicted class members with respect to all class members, whereas specificity  $Sp$  is the proportion of correctly predicted non-class members with respect to all non-class members.

$$Se = |\hat{C} \cap C| / |\hat{C}|,$$

$$Sp = |\bar{\hat{C}} \cap \bar{C}| / |\bar{\hat{C}}|,$$

where  $\bar{\phantom{x}}$  denotes the set complement with respect to  $O$ . For a good classification result, sensitivity and

specificity both should be high. The function

$$g = 1 - (1 - Se)^2 - (1 - Sp)^2$$

was therefore chosen as optimisation criterion.

#### 64 **Generalised ROC-Analysis**

65 For the optimisation of cutoffs in case of more than one prototype, the idea of ROC-analysis was  
 66 generalised. All but one of the prototypes with their corresponding cutoffs have been fixed, and the  
 67 remaining prototype and its cutoff are optimised in each step, only. Without loss of generality, let the  
 68 prototypes  $x_i$  and their cutoffs  $f_i$  be fixed for  $i = 1, \dots, k - 1$ . For a prototype  $x_k$ , the cutoff  $f_k$  can be  
 69 optimised by generalising the ROC principle.

According to the model

$$\hat{C} = \bigcup_{i=1}^k \hat{C}_i = \bigcup_{i=1}^{k-1} \hat{C}_i \cup \hat{C}_k,$$

for this particular case, sensitivity and specificity can be written as

$$Se = |\hat{C} \cap C| / |C| = \left( \left| \bigcup_{i=1}^{k-1} \hat{C}_i \cap C \right| + \left| \hat{C}_k \cap \overline{\bigcup_{i=1}^{k-1} \hat{C}_i} \cap C \right| \right) / |C|, \text{ and}$$

$$Sp = |\overline{\hat{C}} \cap \overline{C}| / |\overline{C}| = \left( \left| \overline{\bigcup_{i=1}^{k-1} \hat{C}_i} \cap \overline{C} \right| - \left| \hat{C}_k \cap \overline{\bigcup_{i=1}^{k-1} \hat{C}_i} \cap \overline{C} \right| \right) / |\overline{C}|.$$

70 From this partition it becomes obvious that the analysis needs to be run over all remaining points, not  
 71 predicted as class members according to the first  $k - 1$  prototypes, only.

#### 72 **Discrete optimisation**

The estimation of the parameters  $\{(x_i \in C, f_i \in \mathbb{R} \geq 0), i = 1, \dots, k\}$  of the given model based on training  
 data  $O = \{o_j \in O, j = 1, \dots, n\}$  with  $|C| = m$  leads to a discrete optimisation problem: Choose a subset of  
 $k$  objects out of the  $m$  training class objects as prototypes, and  $k$  corresponding cutoffs out of the actual  
 distances between training objects and prototypes which jointly optimise the objective function  $g$ :

$$\max_{(x_i, f_i)} (1 - (1 - Se)^2 - (1 - Sp)^2).$$

However, the direct solution of the entire optimisation task would take a prohibitive amount of  
 computation time and memory consumption. There are  $\binom{m}{k}$  possible choices of prototypes and  $n$  potential  
 cutoffs for every prototype. For an effective assessment of the best cutoff, prior sorting of all distances

between the objects and the prototype considered is necessary. This can be done in  $\mathcal{O}(n \cdot \log(n))$  time. Furthermore, the number of prototypes is intended to be small, and for a well-posed problem, the numbers of class members and of non-members both should be of the same order. Presuming  $k \ll n$  and  $m, n-m \in \mathcal{O}(n)$  therefore, results in a time consumption of

$$\binom{m}{k} \cdot n \cdot \log(n) \cdot k \in \mathcal{O}(n^{k+1} \cdot \log(n)).$$

73 The exact solution of this problem needs exponential time in  $k$  and polynomial time with a high  
74 exponent with respect to  $n$ . The below presented greedy algorithm, however, reduces time consumption in  
75 both aspects.

## 76 **2.2 Algorithm and implementation**

### 77 ***Greedy algorithm***

78 The complete solution of the multi-dimensional discrete optimisation problem would soon exceed  
79 available computer resources. For a reasonable determination of sufficiently good  $k$  prototypes and their  
80 corresponding cutoffs, a greedy algorithm (Cormen et al., 2009) was developed. The problem was reduced,  
81 therefore, to repeatedly choosing only one of the  $k$  prototypes anew. An approximate solution of the  
82 entire problem can be found this way. The main principles of the algorithm are here described in the  
83 first instance. Several improvements and practical simplifications will be incorporated in a subsequent  
84 paragraph. Starting with a random choice of  $k$  objects of the training set as preliminary prototypes, and  
85 with cutoffs all initialised with  $-1$ , an iteration process is started. A negative radius of a hypersphere  
86 corresponds to the empty set. At the start, there are no predicted class members.

87 In each cycle, it is investigated for all prototypes whether a better cutoff can be found, or there is  
88 one of the residual objects which improves  $g$  after substitution and appropriate cutoff-estimation by the  
89 generalised ROC-procedure. The substitution is manifested for the prototype for which the maximum  
90 improvement of the ROC-criterion  $g$  is reached, if possible. Thus, the criterion-value  $g$  is actually  
91 improved in each step.

92 If no better  $g$ -value is found the algorithm stops. The estimation of the classification parameters  
93 has finished therewith. The main ideas of the greedy algorithm are summarised in the pseudo code (see  
94 Tab. 1).

### 95 ***Runtime analysis***

96 The algorithm considers substituting one out of  $k$  prototypes and its associated cutoff in each iteration-cycle.  
97 All  $m - k$  remaining class members and the one, old respective prototype are potential new prototypes,  
98 and the distances from the current prototype to all  $n$  objects of the training data might be cutoffs. Again,

**Table 1.** Pseudocode of the greedy algorithm.

---

Start      random choice of preliminary prototypes with negative cutoffs

Iteration    repeat

            for  $i=1$  to  $k$  do

            { fix all other prototypes and cutoffs, consider prototype  $i$

            { for all remaining class members and prototype  $i$  itself

            { substitute prototype  $i$

            { optimise cutoff by generalised ROC-Analysis

            { calculate criterion

            if any improvement then replace by best new prototype and cutoff

            else stop

            end

---

99 cutoff assessment needs  $\mathcal{O}(n \cdot \log(n))$  time because of prior sorting. Furthermore, numerous applications  
100 of the algorithm to example data support the assumption that the number of iteration steps  $j$  is of the same  
101 order as the number of prototypes. This will also be exemplified by data in Sect. 4.

Presuming  $m, n-m \in \mathcal{O}(n)$  as before, the time consumption of the entire algorithm is thus

$$k \cdot (m - k + 1) \cdot n \cdot \log(n) \cdot j \in \mathcal{O}(k^2 \cdot n^2 \cdot \log(n)).$$

102 Finally, the greedy algorithm needs polynomial time in both parameters. In view of  $k \ll n$ , the practical  
103 overall order is  $\mathcal{O}(n^2 \cdot \log(n))$ .

#### 104 **Practical implementation**

105 The depicted algorithm enables a quick estimation of good prototypes and respective cutoffs. However,  
106 some steps of the general algorithm are redundant and should be skipped to make it faster. Other segments  
107 may be simplified or rearranged. The following four aspects below were treated. Thus, the final algorithm  
108 works more effectively. The computational effort could be reduced in a practically pronounced extent this  
109 way, although it does not change its asymptotic order.

- 110 (1) The first idea of enhancing the performance of the algorithm is a triangle-like successive extension  
111 of the set of prototypes at the start. In the first  $k$  iteration cycles, only one new prototype is added.  
112 The treatment of more than one dummy is not necessary. Given the complete distance matrix of  
113 the training set, a standard ROC analysis is performed for every class-member in the first cycle.  
114 The object achieving the highest ROC-criterion is then chosen as starting prototype. In each of the  
115 following  $k - 1$  iteration cycles, one dummy only, i.e. a randomly chosen prototype with negative  
116 cutoff, is added until all  $k$  prototypes are gained.

- 117 (2) The prototypes are permuted after each iteration step to store the recently renewed object with the  
 118 highest index  $j \leq k$  according to (1). Thus, at most only  $j - 1$  preliminary prototypes need to be  
 119 inquired for better ones in the following cycle. Prototype  $j$  has proven to be the optimum choice in  
 120 the previous cycle.
- 121 (3) When performing a generalised ROC-analysis for the prototype  $i$ , only the distances to those objects  
 122 need to be examined which are not predicted as class members according to any of the other just  
 123 fixed prototypes (see Sect. 2.1).
- 124 (4) In some occasional instances, the algorithm might even stop when some of the dummies have  
 125 not been substituted. The best possible classification may already have been reached with fewer  
 126 prototypes than envisaged beforehand. If this has happened, all remaining dummies are removed  
 127 from the final parameter set.

### 128 ***Application to multi class problems***

129 The classification algorithm described so far accomplishes the selection of prototypes and cutoffs for one  
 130 pre-agreed class. In order to solve a multi class problem, this procedure is performed separately for every  
 131 single class in a one against all approach. Thus, at the end of the analysis of training data, only a few  
 132 prototypes and respective cutoffs have been selected for each of the classes.

In anticipation of subsequent multi class diagnostics, the empirical functions of positive predictive values (PPV) are supplementally stored for each of the prototypes. The positive predictive values can easily be derived in the course of ROC-analysis. It is the proportion of correctly predicted class members with respect to all predicted class members, both justified by the prototype under consideration. For the prototype  $x_i$  of class  $C_j$  at distance  $d^*$ , it is

$$PPV_i = PPV_i(d^*) = \frac{|\{y \in C_j : d(y, x_i) \leq d^*\}|}{|\{y \in O : d(y, x_i) \leq d^*\}|}.$$

133 Given the estimated parameters of all single class problems and the corresponding empirical PPV  
 134 functions, the distances of any object of unknown class to all prototypes must be calculated first. Only  
 135 those prototypes with distances not higher than their corresponding cutoffs are further considered. If  
 136 exactly one prototype is left, its class is predicted for the object. If more than one prototype comes into  
 137 question, the object is then assigned to the class of the prototype with the highest related PPV. If all  
 138 distances were higher than the corresponding cutoffs, the object is assigned to an additional residual class.  
 139 This result should clarify that a reliable classification of that object is not possible based on the data given.

140 The ruling based on PPV facilitates a balanced diagnostic decision for all classes; and the extraneous  
 141 residual class ensures a high specificity of classification even if it is tried for an object that indeed does

142 not belong to any of the trained classes.

### 143 **Software**

144 The elucidated multi prototype classifier was implemented in R (R Core Team, 2012), version 2.15.2  
145 (2012-10-26). The sources of the corresponding R-functions “fit.prototypes.r” and “predict.prototypes.r”  
146 will be made available on request.

## 147 **3 EVALUATION EXAMPLES**

148 The performance of the classification system presented in this paper is demonstrated at various application  
149 examples, and in comparison with other methods. In this context, time expenditure and accuracy of  
150 prediction are the main aspects of evaluation.

151 For the comparison of final classification results, six multi class problems have been selected. They  
152 are based on example data sets provided by the UCI Machine Learning Repository (Bache and Lichman,  
153 2013). All data sets comprise moderately many features and classes. All features are numeric and have  
154 been treated as continuous data. The number of included objects enables a serious and fast classification.  
155 One data set (Statlog shuttle) comprises many more objects than the others. So, distinctions in computation  
156 time requirements can be illustrated better. The chosen data examples cover various fields of application:  
157 biomedical data, image data, and physical data.

### 158 **3.1 Data sets**

159 *Cardiotocography data set (Ayres-de Campos et al., 2000)*. The data set consists of 21 measured features  
160 of fetal heart rate and uterine contraction on 2126 cardiotocograms (“Cardiotocography”) classified by  
161 expert obstetricians with respect to both, a morphologic pattern and a fetal state. In this collection fetal  
162 state classes (NSP) are considered, only. Their frequencies are

1	2	3	.
1655	295	176	.

164 *Statlog image segmentation data set (Brodley, 1990)*. Segments from seven images (classes) were  
165 evaluated for several pixel measures. Each of the 2310 instances represents a 3x3 region. Attribute no. 3  
166 was removed because of identical entries. So, this data set (“Statlog Image Segmentation”) comprises 18  
167 features to predict seven equally sized classes:

1	2	3	4	5	6	7	.
330	330	330	330	330	330	330	.

169 *Statlog shuttle data set (Catlett, 1994)*. NASA Shuttle database deals with the positioning of radiators  
170 in the Space Shuttle. In the original, it consists of 43500 training objects and 14500 test objects. Each  
171 instance is described by nine continuous attributes and is assigned to one of seven classes with the  
172 frequencies

173           1       2    3    4       5       6    7  
           45586  50  171  8903  3267  10  13

174 Classes 6 and 7 are too small for substantial classification and are omitted afterwards which consequently  
 175 results in the frequencies

176           1       2    3    4       5  
           45586  50  171  8903  3267

177 All 57977 remaining objects are united into one data set (“Statlog Shuttle”) in advance of subsequent  
 178 cross-validation.

179       *Statlog vehicle silhouettes data set (Siebert, 1987)*. The data set (“Statlog Vehicle Silhouettes”)  
 180 provides 18 features extracted from the silhouettes of four types of vehicle. The vehicle may have been  
 181 viewed from one of many different angles. The overall number of examples amounts to 846, its class  
 182 frequencies are:

183           bus  opel  saab  van  
           218  212  217  199

184       *Waveform database generator (version 1) data set (Breiman et al., 1984)*. This artificial data set  
 185 (“Waveform Data Generator”) was simulated for classification purpose (for details see Breiman et al.  
 186 (1984)). Three classes of waves were generated from a combination of two of three “base” waves. Each  
 187 instance consists of 21 attributes, all of which include noise. A total of 5000 instances was uniformly  
 188 partitioned into the classes:

189           1       2       3  
           1657  1647  1696

190       *Yeast data set (Horton and Nakai, 1996, 1997) (modified)*. This data set was collected for predicting  
 191 the cellular localisation sites of proteins. In the original data, eight numeric attributes of altogether 1484  
 192 objects were provided, and proteins were divided into ten classes (Horton and Nakai, 1996). The numbers  
 193 of objects of several classes, however, were fairly low. For the purpose under investigation, all three  
 194 membrane proteins (ME1, ME2, ME3) have been united to class MEM, therefore. Furthermore, all  
 195 other small groups (EXC, VAC, POX, ERL) have been pooled to a residual class RES. The original class  
 196 frequencies were

197           CYT  MIT  NUC  ME1  ME2  ME3  ERL  EXC  POX  VAC  
           463  244  429  44  51  163  5  35  20  30

198 whereas the frequencies of the modified data set (“Yeast Protein Localisation”) are

199           CYT  MIT  NUC  MEM  RES  
           463  244  429  258  90

201 Descriptive parameters of all applied data sets are gathered in Tab. 2.

**Table 2.** Numbers of objects, features, and classes of evaluation data sets.

Data Set	Number of		
	Objects	Features	Classes
Cardiotocography	2126	21	3
Statlog Image Segmentation	2310	18	7
Statlog Shuttle	57977	9	5
Statlog Vehicle Silhouettes	846	18	4
Waveform Data Generator	5000	21	3
Yeast Protein Localisation	1484	8	5

### 202 **3.2 Methods**

203 For demonstration and evaluation, the described classification algorithm is compared with several widely  
204 used classifiers. Data-based and feature-based methods have been selected as well to get a general  
205 impression of the performance.

206 *Multi prototype classifier.* Multi prototype classification (“PRO”) has been performed for all data sets,  
207 each with a series of different parameters. Euclidean, Manhattan, or Canberra distances were applied to  
208 the data. The maximum number of prototypes was set to 3, 5, or 7.

209 *k-Nearest neighbour classifier (Tan et al., 2005).* The implementation of the k-nearest neighbour  
210 algorithm (“kNN”) of the R package “knnflex” (Brooks, 2007) was utilised. Again, Euclidean, Manhattan,  
211 or Canberra distances were applied to the data. The number of neighbours considered for classification  
212 was chosen as 1, 3, or 5.

213 *Support vector machine (Kohonen, 2001).* The support vector machine algorithm (“SVM”) of the R  
214 package “e1071” (Meyer et al., 2012) was utilised.

215 *CART (Breiman et al., 1984).* This classical variant of decision trees (“CART”) provided by the R  
216 package “rpart”(Therneau et al., 2012) was applied.

217 *C4.5 (Quinlan, 1993).* The implementation of the R package “RWeka” (Hornik et al., 2009; Witten  
218 and Frank, 2005) was used for this algorithm (“C4.5”).

### 219 **3.3 Validation**

220 *Limited 10-fold cross-validation.* In each of the data analyses, a limited 10-fold cross-validation has  
221 been performed. The data set is randomly partitioned into ten almost equal-sized subsets, therefore. In  
222 each of the ten steps, one subset serves as test data as in usual cross-validation (Tan et al., 2005; Hastie  
223 et al., 2009). The remaining nine subsets together constitute the preliminary training data. It is randomly  
224 reduced to maximum 100 objects of each class if more are present. With this limitation of the training  
225 data sets, the estimation of model parameters is based on the information of at most 100 objects of each  
226 class for all classification procedures. Although it might not lead to the best classification result possible,

**Table 3.** Optimum distance measures and numbers of prototypes or neighbours, respectively.

Data Set	Multi prototypes		k-Nearest neighbours	
	Metric	No. of prototypes	Metric	No. of neighbours
Cardiotocography	Canberra	7	Canberra	1
Statlog Image Segmentation	Manhattan	7	Manhattan	1
Statlog Shuttle	Euklidean	7	Manhattan	1
Statlog Vehicle Silhouettes	Euklidean	7	Manhattan	1
Waveform Data Generator	Euklidean	3	Manhattan	5
Yeast Protein Localisation	Manhattan	5	Manhattan	5

227 it should do for comparison.

228 *Adjusted Rand index.* The summarised true and predicted class frequencies are compared by the  
229 adjusted Rand index (“ARI”) (Vinh et al., 2010; Hubert and Arabie, 1985; Vinh et al., 2009). It is corrected  
230 for the bias of different class sizes and is therefore considered as main criterion of success.

231 *Accuracy.* The plain accuracy (“ACC”) or naive percentage as the proportion of correctly classified  
232 objects disregards different numbers of objects in the various classes. This index is calculated in percentage  
233 terms for additional consideration.

234 *Number of iteration cycles.* In case of multi prototype classifier application, the number of iteration  
235 cycles is recorded for every single class approach. This oblique processing parameter affects the effective  
236 calculation time.

237 *Computation time.* Computation time is recorded for training and prediction parts of the calculations,  
238 separately. It should demonstrate the operability of the applied methods for sundry applications.

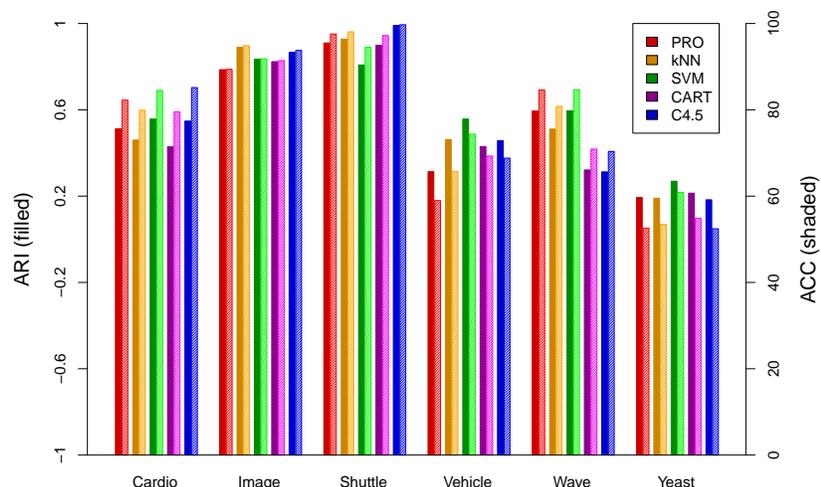
239 *Software.* All calculations were performed using R (R Core Team, 2012), Version 2.15.2 (2012-10-26).

## 240 **4 RESULTS**

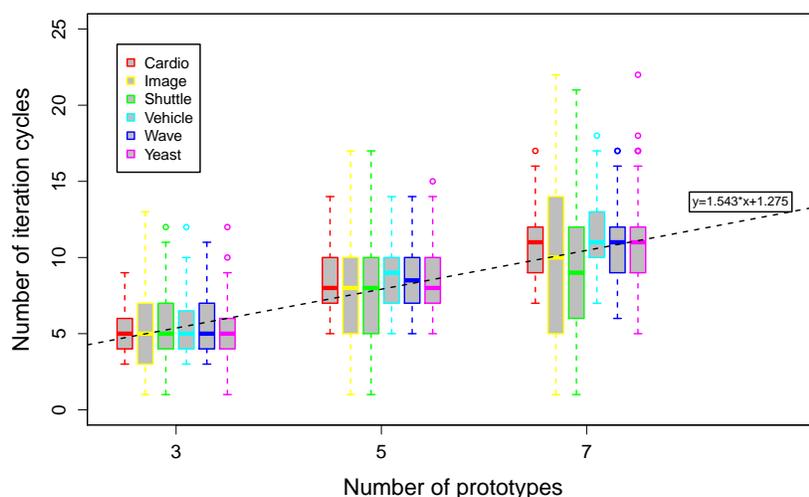
241 All data sets have been analysed as described in Sect. 3.2. In case of multi prototype or k-nearest  
242 neighbour classifiers, the best classifications with respect to distance measure and number of prototypes  
243 or neighbours, respectively, were used for comparison. Adjusted Rand index was applied as selection  
244 criterion. The best choices of the processing parameters are listed in Tab. 3.

245 The classification performances of the compared methods were measured by ARI and ACC. All  
246 results are visualised in Fig. 1. They vary with data sets and with methods. The values of ARI range from  
247 0.183 .. 0.269 for the yeast data to 0.808 .. 0.991 for the shuttle data. The corresponding ACC ranges are  
248 52.4 .. 60.8 (yeast) and 94.5 .. 99.7 (shuttle).

249 In every performed basic multi prototype classification, i.e. single class analysis, the number of  
250 iteration cycles was recorded. Fig. 2 shows boxplots for all data sets dependent on the number of  
251 prototypes each class. According to the applied cross-validation procedure, one boxplot summarises ten



**Figure 1.** Adjusted Rand indices and accuracies of the classification results.



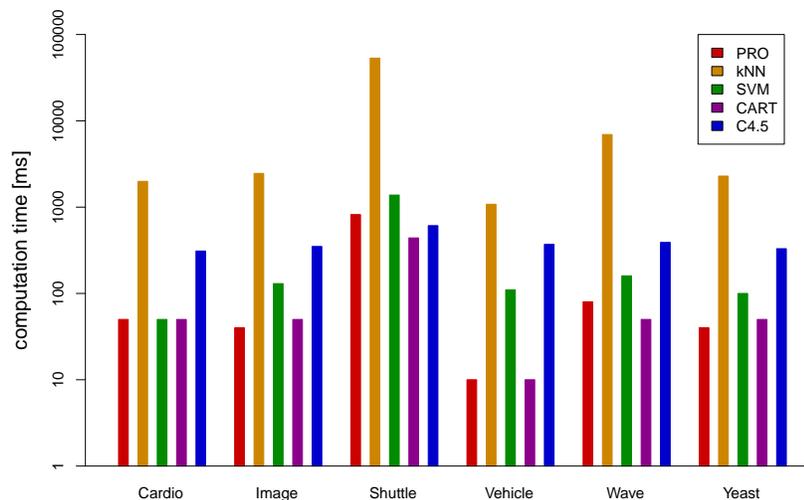
**Figure 2.** Number of iteration cycles of the multi prototype classifier.

252 repetitions of the number of cycles for each of the classes of the underlying data set, and for each of the  
 253 metrics investigated. The estimated parameters of a linear regression of the entire resulting data are 1.543  
 254 (inclination) and 1.275 (intercept). The data support the assumption that the number of iteration cycles is  
 255 of the same order as the processed number of prototypes, although it cannot be evidenced with it.

256 The computational effort for the prediction part of classification is of paramount interest in diagnostic  
 257 applications. There are large speed differences between the methods. The corresponding results are  
 258 depicted in Fig. 3 on a logarithmic time scale, therefore.

## 259 5 DISCUSSION

260 In this paper, a novel, universal distance-based classification procedure has been introduced. For each  
 261 class to be recognised, this classifier computes a number of prototypes and corresponding cutoffs from



**Figure 3.** Computation time need of the prediction part.

262 the training data. The approximate solution of the underlying discrete optimisation problem has been  
 263 realised by implementing a greedy algorithm which runs in  $\mathcal{O}(k^2 n^2 \log(n))$  time where  $k$  is the number  
 264 of prototypes to be selected and  $n$  is the number of training objects.

265 In contrast to the concepts of centroids or medoids, it is not necessary that the centres of the hyper-  
 266 spheres coincide with regions of high point density. The multi prototype classifier models the shape of the  
 267 points belonging to the class by an overlap of different-sized hyperspheres. This global distance-based  
 268 approach is universally applicable. It is based on a simple geometric model. No other assumptions, e.g.  
 269 about the distribution of point distances or about clusters or subclusters, were made. Nevertheless, the  
 270 choice of an appropriate distance measure is of great importance for the classification performance, as  
 271 always in distance-based methods.

272 The presented method can be utilised in case of multi-class problems as well. For this purpose, a one  
 273 against all approach has been performed. This makes a future extension of the entire diagnostic method  
 274 for additional classes easily possible. An object has been finally assessed to that class of maximum  
 275 positive predictive value when in doubt. Objects not recognised as a member of any of the classes have  
 276 been assigned to an additional residual class.

277 For classification methods applied in diagnostics, the introduction of a residual class opens the chance  
 278 to avoid unnecessary misdiagnoses and to examine those objects by other tools. This ensures a high level  
 279 of confidence. Although the extensive evaluation of more pairwise classification results may enhance  
 280 the overall performance in general (Allwein et al., 2000), the choice of a one against all approach seems  
 281 to be sufficient. The results are comparable to those of other classification procedures and demonstrate  
 282 excellent agreements of the class labels. Thus, multi prototype classifier can be used alternatively to other  
 283 techniques.

284

285 Classification results of some few methods applied to various data sets have been compared by the  
286 poor accuracy index, and by the adjusted Rand index. ACC is biased by differently sized classes. In  
287 general, ARI is a good and reliable measure of classification success and should be preferred. In the cases  
288 considered in this paper, however, both indices would not lead to dissent conclusions when used apartly.

289 The classification performances visualised in Fig. 1 vary with both, data sets and methods. So,  
290 the decision which procedure should be applied depends on the specific task. The results convey the  
291 impression that multi prototype classifier is comparable to the tools considered and works more or less  
292 effectively, like others else.

293

294 Time consumption is generally hard to compare between different computer programs. Each of them  
295 is possibly programmed by other persons, and in different contexts, and may differ in effectivity. This  
296 applies to the specific implementations of the classifiers described in Sect. 3.2, too. Furthermore, PRO  
297 is interpreting R-code, only. The other examined methods make partly use of compiled segments. The  
298 differences in computation time, as presented in Fig. 3, should consequently be interpreted with care.  
299 However, some patterns are expected and plausible.

300 For all data sets, kNN is the slowest prediction method. The distances of an investigated object to  
301 all objects of the reference set must be calculated for it. Just for that reason, PRO was developed. The  
302 prediction using C4.5, with the exception of shuttle data, was also relatively slow. This might be due to  
303 relatively large trees necessary for good classification. PRO, however, has proven to be one of the fastest  
304 methods. This is an important benefit for potential utilization such as in high-throughput diagnostics.

305

306 The findings of this exploration can be condensed in three main statements. Multi prototype classifier  
307 is a universal classification algorithm which is comparable to the other examined procedures concerning its  
308 overall performance. With the introduction of a residual class, PRO realises the opportunity to tag objects  
309 which cannot be classified at a sufficient level of confidence. The depicted multi prototype classification  
310 technique enables a rapid class prediction even for large amounts of data.

## 311 **ACKNOWLEDGMENTS**

312 The author would like to acknowledge UCI Machine Learning Repository (Bache and Lichman, 2013).  
313 These data sets are very helpful for developing, testing, and comparing algorithms for data analysis.  
314 Providing that material will also help readers to practically reproduce presented results.

## 315 REFERENCES

- 316 Allwein, E. L., Schapire, R. E., and Singer, Y. (2000). Reducing multiclass to binary: a unifying approach  
317 for margin classifiers. *Journal of Machine Learning Research*, 1:113–141.
- 318 Ayres-de Campos, D., Bernardes, J., Garrido, A., Marques-de Sá, J., and Pereira-Leite, L. (2000). SisPorto  
319 2.0: a program for automated analysis of cardiocograms. *J Matern Fetal Med*, 9(5):311–8.
- 320 Bache, K. and Lichman, M. (2013). UCI machine learning repository.
- 321 Breiman, L., Friedman, J., Stone, C., and Olshen, R. (1984). *Classification and Regression Trees*.  
322 Chapman & Hall, New York.
- 323 Brodley, C. (1990). Statlog (Image Segmentation) Data Set. Vision Group, University of Massachusetts.
- 324 Brooks, A. D. (2007). *knflex: A more flexible KNN*. R package version 1.1.1.
- 325 Catlett, J. (1994). Statlog (Shuttle) Data Set. Basser Department of Computer Science, University of  
326 Sydney, N.S.W., Australia.
- 327 Cormen, T. H., Leiserson, C. E., Rivest, R. L., and Stein, C. (2009). *Introduction to Algorithms*. The MIT  
328 Press, Cambridge, MA, USA, 3rd edition.
- 329 De Bruyne, K., Slabbinck, B., Waegeman, W., Vauterin, P., De Baets, B., and Vandamme, P. (2011).  
330 Bacterial species identification from MALDI-TOF mass spectra through data analysis and machine  
331 learning. *Syst Appl Microbiol*, 34(1):20–9.
- 332 Gordon, A. (1999). *Classification*. Monographs on Statistics and Applied Probability. Chapman &  
333 Hall/CRC.
- 334 Greiner, M., Pfeiffer, D., and Smith, R. (2000). Principles and practical application of the receiver-  
335 operating characteristic analysis for diagnostic tests. *Preventive Veterinary Medicine*, 45(1-2):23–41.
- 336 Hanley, J. A. and McNeil, B. J. (1982). The meaning and use of the area under a receiver operating  
337 characteristic (roc) curve. *Radiology*, 143(1):29–36.
- 338 Hastie, T., Tibshirani, R., and Friedman, J. (2009). *The elements of statistical learning: data mining,*  
339 *inference and prediction*. Springer, 2 edition.
- 340 Hong, S. P., Shin, S., Lee, E. H., Kim, E. O., Ji, S. I., Chung, H. J., Park, S. N., Yoo, W., Folk, W. R., and  
341 Kim, S. (2008). High-resolution human papillomavirus genotyping by MALDI-TOF mass spectrometry.  
342 *Nat. Protocols*, 3(9):1476–1484.
- 343 Hornik, K., Buchta, C., and Zeileis, A. (2009). Open-source machine learning: R meets weka. *Comput.*  
344 *Stat.*, 24(2):225–232.
- 345 Horton, P. and Nakai, K. (1996). A probabilistic classification system for predicting the cellular localization  
346 sites of proteins. In *Proceedings of the Fourth International Conference on Intelligent Systems for*  
347 *Molecular Biology*, pages 109–115. AAAI Press.

- 348 Horton, P. and Nakai, K. (1997). Better prediction of protein cellular localization sites with the k nearest  
349 neighbors classifier. In *Proceedings of the Fifth International Conference on Intelligent Systems for*  
350 *Molecular Biology*, pages 147–152. AAAI Press.
- 351 Hubert, L. and Arabie, P. (1985). Comparing partitions. *Journal of Classification*, 2(1):193–218.
- 352 Karger, A., Ziller, M., Bettin, B., Mintel, B., Schares, S., and Geue, L. (2011). Determination of  
353 serotypes of shiga toxin-producing escherichia coli isolates by intact cell matrix-assisted laser desorption  
354 ionization-time of flight mass spectrometry. *Applied and Environmental Microbiology*, 77(3):896–905.
- 355 Kohonen, T. (2001). *Self-Organizing Maps*. Springer, Berlin, 3rd edition.
- 356 Martinetz, T. M. and Schulten, K. J. (1991). A "neural-gas" network learns topologies. In Kohonen,  
357 T., M"akisara, K., Simula, O., and Kangas, J., editors, *Artificial Neural Networks*, pages 397–402.  
358 North-Holland, Amsterdam, North-Holland, Amsterdam.
- 359 Meyer, D., Dimitriadou, E., Hornik, K., Weingessel, A., and Leisch, F. (2012). *e1071: Misc Functions of*  
360 *the Department of Statistics (e1071), TU Wien*. R package version 1.6-1.
- 361 Quinlan, J. R. (1993). *C4.5: programs for machine learning*. Morgan Kaufmann Publishers Inc., San  
362 Francisco, CA, USA.
- 363 R Core Team (2012). *R: A Language and Environment for Statistical Computing*. R Core Team, Vienna,  
364 Austria. ISBN 3-900051-07-0.
- 365 Siebert, J. P. (1987). Vehicle recognition using rule based methods. Turing Institute Research Memorandum  
366 TIRM-87-018. Technical report, Turing Institute.
- 367 Tan, P., Steinbach, M., and Kumar, V. (2005). *Introduction to Data Mining, (First Edition)*. Addison-  
368 Wesley Longman Publishing Co., Inc., Boston, MA, USA.
- 369 Therneau, T., Atkinson, B., and Ripley, B. (2012). *rpart: Recursive Partitioning*. R package version 4.1-0.
- 370 Vinh, N. X., Epps, J., and Bailey, J. (2009). Information theoretic measures for clusterings comparison: is  
371 a correction for chance necessary? In *Proceedings of the 26th Annual International Conference on*  
372 *Machine Learning, ICML '09*, pages 1073–1080, New York, NY, USA. ACM.
- 373 Vinh, N. X., Epps, J., and Bailey, J. (2010). Information theoretic measures for clusterings comparison:  
374 Variants, properties, normalization and correction for chance. *J. Mach. Learn. Res.*, 11:2837–2854.
- 375 Witten, I. H. and Frank, E. (2005). *Data Mining: Practical Machine Learning Tools and Techniques, Sec-*  
376 *ond Edition (Morgan Kaufmann Series in Data Management Systems)*. Morgan Kaufmann Publishers  
377 Inc., San Francisco, CA, USA.
- 378 Wong, J., Schwahn, A., and Downard, K. (2010). FluTyper - an algorithm for automated typing and  
379 subtyping of the influenza virus from high resolution mass spectral data. *BMC Bioinformatics*,  
380 11(1):266.

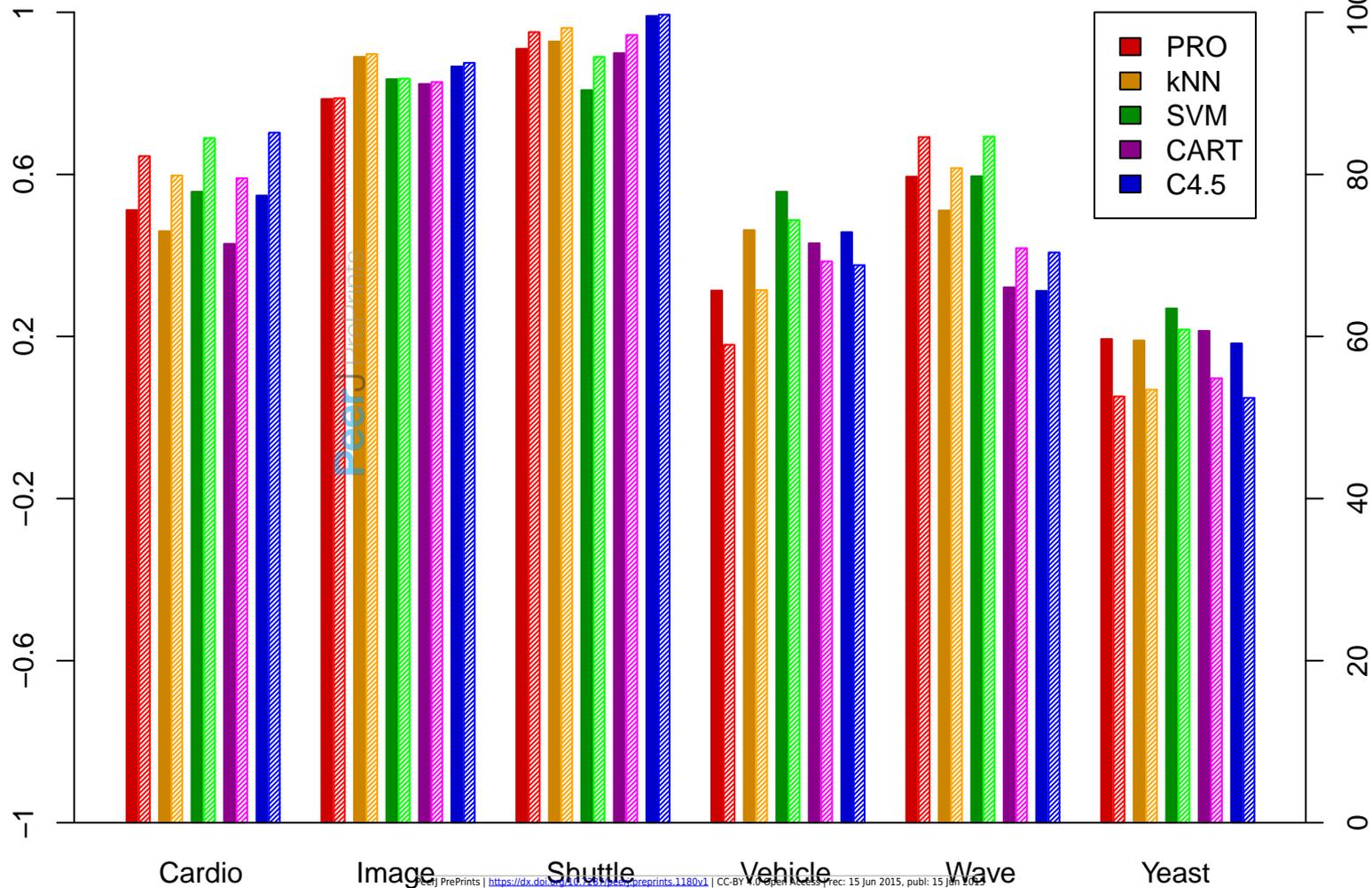
381 Ziller, M., Gussmann, M., and Karger, A. (2011). Mathematical aspects of classifying mass spectra.  
382 diagnostic discrimination of influenza-subtypes. In Dress, A., Biebler, K.-E., Cieslik, D., and Spillner,  
383 A., editors, *The Math of Flu*, volume 17 of *Biometrie und Medizinische Informatik - Greifswalder*  
384 *Seminarberichte*, pages 101–123. Shaker, Aachen.

**Figure 1** (on next page)

Figure 1

Adjusted Rand indices and accuracies of the classification results.

ARI (filled)



ACC (shaded)

**Figure 2** (on next page)

Figure 2

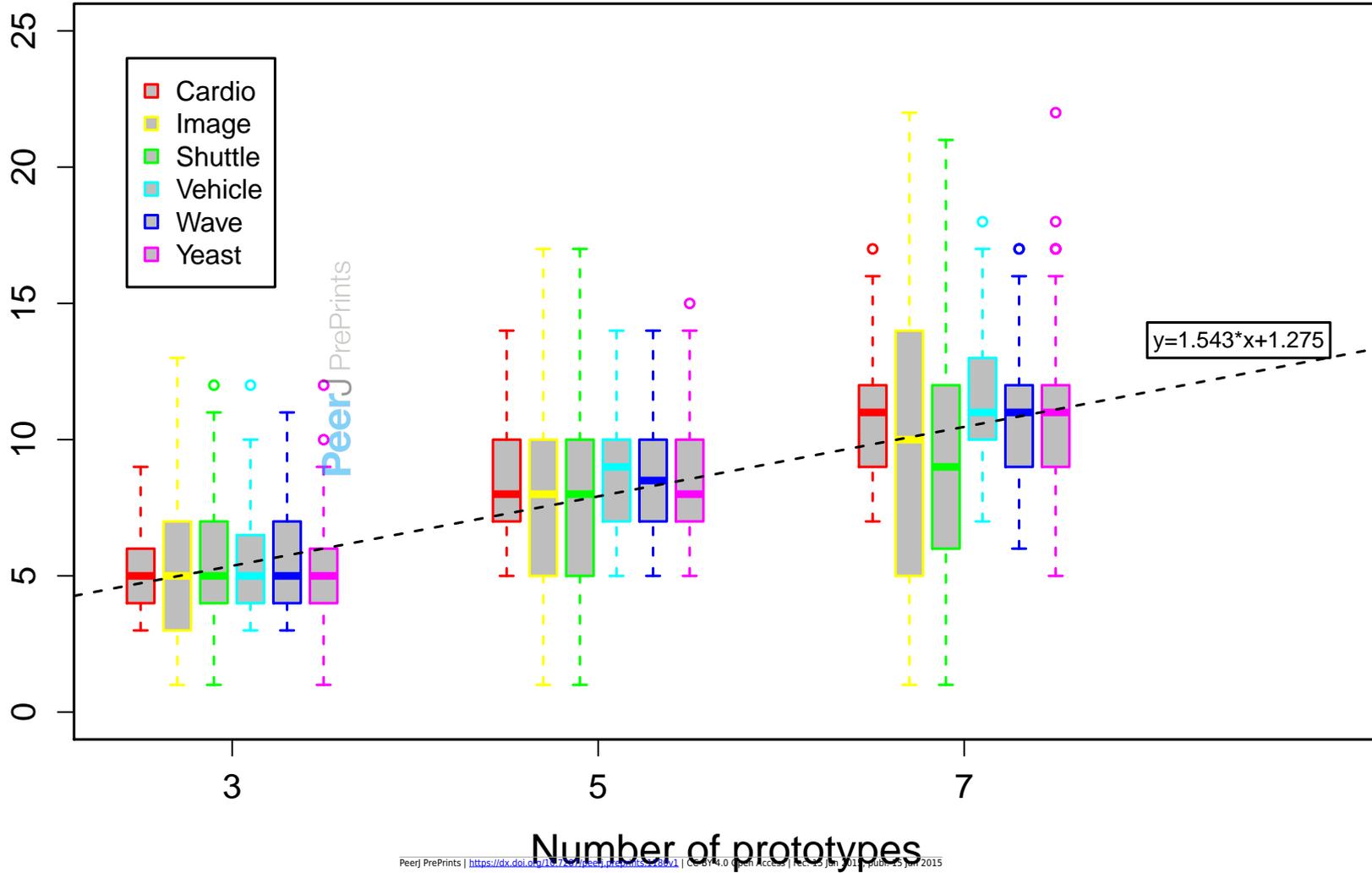
Number of iteration cycles of the multi prototype classifier.

Number of iteration cycles

- Cardio
- Image
- Shuttle
- Vehicle
- Wave
- Yeast

PeerJ PrePrints

$$y=1.543*x+1.275$$



**Figure 3** (on next page)

Figure 3

Computation time need of the prediction part.

