

# Experimental study on prototype optimisation algorithms for prototype-based classification in vector spaces

M. Lozano<sup>a,\*</sup>, J.M. Sotoca<sup>a</sup>, J.S. Sánchez<sup>a</sup>, F. Pla<sup>a</sup>, E. Pełalska<sup>b,c</sup>, R.P.W. Duin<sup>b</sup>

<sup>a</sup>Dept. Lenguajes y Sistemas Informáticos, Universitat Jaume I, Campus Riu Sec, 12071 Castellón, Spain

<sup>b</sup>Faculty of Electrical Engineering, Mathematics and Computer Science, Delft University of Technology, Mekelweg 4, 2628 CD Delft, The Netherlands

<sup>c</sup>School of Computer Science, The University of Manchester, Manchester M13 9PL, UK

Received 13 February 2006; accepted 6 April 2006

## Abstract

Prototype-based classification relies on the distances between the examples to be classified and carefully chosen prototypes. A small set of prototypes is of interest to keep the computational complexity low, while maintaining high classification accuracy. An experimental study of some old and new prototype optimisation techniques is presented, in which the prototypes are either selected or generated from the given data. These condensing techniques are evaluated on real data, represented in vector spaces, by comparing their resulting reduction rates and classification performance.

Usually the determination of prototypes is studied in relation with the nearest neighbour rule. We will show that the use of more general dissimilarity-based classifiers can be more beneficial. An important point in our study is that the adaptive condensing schemes here discussed allow the user to choose the number of prototypes freely according to the needs. If such techniques are combined with linear dissimilarity-based classifiers, they provide the best trade-off of small condensed sets and high classification accuracy.

© 2006 Pattern Recognition Society. Published by Elsevier Ltd. All rights reserved.

*Keywords:* Dissimilarity representation; Prototype selection; Adaptive condensing; EM algorithm; Normal density based classifier; Nearest neighbour rule

## 1. Introduction

An intuitive way of determining the class of an unknown object is by analysing its similarity to a set of prototypes either selected or generated from a given training set (TS) of objects with known class labels. In general, similarities or dissimilarities can be computed either from the raw object observations or based on an intermediate feature representation. A small set of prototypes has the advantage of a low computational cost and small storage requirements, while leading to similar, or even improved, classification performance. Various ways of designing a prototype set can be studied in Euclidean vector spaces. Two families of such optimisation procedures are editing and condensing.

*Editing* is the step in a learning process in charge of increasing the accuracy of predictions, when there is substantial noise in the training data. A basic editing algorithm removes noisy instances, as well as close border cases, eliminating a possible overlap between the regions from different classes and leaving smoother decision boundaries. Wilson introduced the first editing method [1]. Briefly, the  $k$ -nearest neighbour ( $k$ -NN) rule is used to estimate the class of each prototype in the TS and to remove those whose class labels do not agree with the ones judged by the  $k$ -NN rule. This algorithm tries to eliminate mislabelled objects from the TS as well as those near to the decision boundaries. Many researchers have addressed the problem of editing by proposing alternative schemes [2–5].

The *condensing* step aims at selecting a small subset of prototypes without a significant degradation in the classification accuracy. Two main groups of condensing techniques can be distinguished. These are the selective

\* Corresponding author. Tel.: +34 9767 62109; fax: +34 9767 61914.

E-mail addresses: [lozano@uji.es](mailto:lozano@uji.es) (M. Lozano), [sotoca@uji.es](mailto:sotoca@uji.es) (J.M. Sotoca), [sanchez@uji.es](mailto:sanchez@uji.es) (J.S. Sánchez), [pla@uji.es](mailto:pla@uji.es) (F. Pla), [e.m.pełalska@tudelft.nl](mailto:e.m.pełalska@tudelft.nl) (E. Pełalska), [r.p.w.duin@ieee.org](mailto:r.p.w.duin@ieee.org) (R.P.W. Duin).

schemes, which merely select a subset of the original training objects [2,6–9] and the adaptive schemes which modify them [10–15].

This paper discusses prototype optimisation methods, such as editing and condensing, for feature-based representations of classes in the context of prototype-based classification. Traditionally, the 1-NN rule is used for this purpose. It classifies objects based on the minimum distance to the given prototypes. Here, we show that the prototype sets, optimised to guarantee good performance of the 1-NN rule, lead to a higher classification accuracy, when more general dissimilarity-based classifiers are considered, as recently proposed [16,17]. These are weighted linear or quadratic combinations of the (Euclidean) distances computed between the test objects and the prototype sets found by dedicated condensing algorithms. Linear classifiers are especially of interest, since their computational complexity is comparable to that of the 1-NN rule.

Although the algorithms are run in vector spaces, the distances are by no means restricted to the Euclidean metric. Other distances, such as the  $l_p$ -distances,  $d(\mathbf{x}, \mathbf{y}) = (\sum_i |x_i - y_i|^p)^{1/p}$ ,  $p > 0$  (metric for  $p \geq 1$ ), or inner-product based distances,  $d(\mathbf{x}, \mathbf{y}) = 1 - \mathbf{x}^T \mathbf{y} / \|\mathbf{x}\| \|\mathbf{y}\|$ , can be used. Moreover, the selective techniques can easily be applied to non-vectorial data (such as strings, graphs or shapes) provided that the pairwise dissimilarities are derived. This aspect has already been confirmed in earlier studies, in which numerous embeddings of metric and non-metric dissimilarity data (usually derived from non-vectorial representations) were analysed in Euclidean and pseudo-Euclidean spaces [16–18]. There, the prototypes were often chosen at random [18–21]. Other object selection techniques based on feature selection or clustering approaches were studied and appeared to be good for small prototype sets [16,17]. The conclusion of that work is that more general dissimilarity-based classifiers defined by the dissimilarities to a set of selected prototypes are competitive to the NN rule in the context of non-vectorial dissimilarity data.

In this paper we will confirm and extend this conclusion to vectorial representations, where (additionally to selection approaches) adaptive schemes can naturally be applied to create prototypes. These aspects have not been investigated so far. The reason behind our contribution is that the 1-NN (or  $k$ -NN) rule is often applied in vector spaces<sup>1</sup> for which editing-and-condensing techniques have been widely studied in order to reduce the computational burden while maintaining high accuracy. Consequently, our paper focusses on vectorial representations and the aspects related to the NN and condensed techniques. It is not our goal to present a general investigation into prototype selection techniques for dissimilarity data resulting from non-vectorial representations; in this case, the reader is referred to Refs. [16,17]. Instead, our goal is to study the applicability of condens-

ing techniques, used in vector spaces to optimise the 1-NN rule, as prototype optimisation techniques for building more general dissimilarity-based classifiers.

We will show that linear or quadratic dissimilarity-based classifiers may successfully replace the 1-NN rule, especially when small sets of prototypes are needed. We will focus on Euclidean distances to maintain the connection with the traditionally used distances for selective condensing schemes in vector spaces. We will also compare these with adaptive reduction algorithms. Consequently, we focus on Euclidean distances in vector spaces and the use of condensing schemes determined in favour of the 1-NN rule. Adaptive schemes are especially of interest, as they can only be used when vectorial representations are available. Consequently, as they offer more flexibility, they should be more beneficial in vector spaces than the selective approaches.

Four condensing techniques are studied. Experiments are conducted to compare their ability to reduce the training size, while maintaining the discriminative power of the optimised prototypes. The classification performance is judged by the 1-NN rule and two more general dissimilarity-based classifiers.

The paper is organised as follows: Section 2 briefly reviews a number of condensing techniques which are used in our study. These are *MaxNCN* [14,23], *Reconsistent* [14], *LVQ* [24] and *MixtGauss* [15]. Section 3 briefly describes the framework of the prototype-based classification methods used for the evaluation of the derived condensed sets. In Section 4, the data sets are presented and the experiments are described, providing quantitative results and a further discussion. Finally, the main conclusions are summarised in Section 5.

## 2. Condensing methods to compare

Assume a TS of  $N$  instances,  $X = \{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N\}$ , representing  $J$  classes,  $C = \{c_1, \dots, c_J\}$ . Each instance  $\mathbf{x}_i$  is a point (a vector) in an  $n$ -dimensional feature space  $\mathcal{R}^n$ . A condensed set consists of  $r$ ,  $r \ll n$ , prototypes, which are either selected or generated from the examples of  $X$ . They are determined to represent efficiently the distributions of the classes and be well discriminative, when used to classify the training objects. Their cardinality should be sufficiently small to reduce both the storage and evaluation time.

A condensed set is said to be *consistent* with respect to a TS if the classification error, estimated by assigning all objects from the TS to the classes of their nearest neighbours in the condensed set, is small; see Refs. [3,7]. Therefore, given a set of prototypes representing the class distribution, the classification rate of the TS can be used to measure the consistency of this set. As the consistent condition is maintained, the smaller the number of prototypes in the condensed set, the better the final result is.

Four condensing algorithms are presented below. These are two selective schemes, *MaxNCN* and *Reconsistent*, and two adaptive schemes, *LVQ* and *MixtGauss*.

<sup>1</sup>This holds since the  $k$ -NN rule is known to be a simple and good classifier for large sample sizes thanks to its theoretical properties [22].

## 2.1. MaxNCN

The geometrical distribution of objects in a vector space may be more informative than the distances between them. The so-called *surrounding neighbourhood-based rules* [23] try to improve the traditional nearest neighbour approach by making use of such information, especially in relation to the training objects which are nearby the decision boundaries. This can be achieved by taking into account not only the proximity of objects to a given input sample but also their *symmetrical distribution* around it. As a result, the surrounding neighbours of an object  $p$  should satisfy two complementary conditions:

- (1) Distance criterion: the neighbours should be as close as possible to  $p$ .
- (2) Symmetry criterion: the neighbours should be distributed as homogeneously as possible around  $p$ .

Chaudhuri [25] proposed the concept of a nearest centroid neighbourhood (NCN), which can be viewed as a particular realisation of the surrounding neighbourhood. Let  $p$  be a given object whose  $k$ -NCNs ( $k$  nearest centroid neighbours) should be found in a TS. These  $k$  neighbours can be determined by the following iterative procedure:

- (1) Find the nearest neighbour of  $p$ . Choose it as the first NCN,  $q_1$ .
- (2) Select the  $i$ th NCN,  $q_i$ ,  $i \geq 2$ , from  $X$  such that the centroid defined by  $q_i$  and the previously selected NCNs,  $q_1, \dots, q_{i-1}$  is the closest to  $p$ .

As designed, the NCN search method is incremental and the objects around a given sample have a geometrical distribution that tends to surround it. Moreover, the region of influence of the NCN is larger than that of the NN, in general. This is illustrated in Fig. 1, where the regions defined by the 4-NCN and the 4-NN are shown for the given point  $p$ .

---

### Algorithm 1 MaxNCN

---

```

for  $i = \text{each\_prototype}(TS)$  do
  neighbours_number[ $i$ ] = 0
  neighbour = next_neighbour( $i$ )
  while neighbour.class =  $i$ .class do
    neighbours_vector[ $i$ ] = Id(neighbour)
    neighbours_number[ $i$ ] ++
    neighbour = next_neighbour( $i$ )
  end while
end for
while Max_neighbours() > 0 do
  Eliminate_Neighbours(id_Max_neighbours)
end while

```

---

The *MaxNCN* technique is based on the concept of NCN and relies on the NCN search algorithm, as presented above.

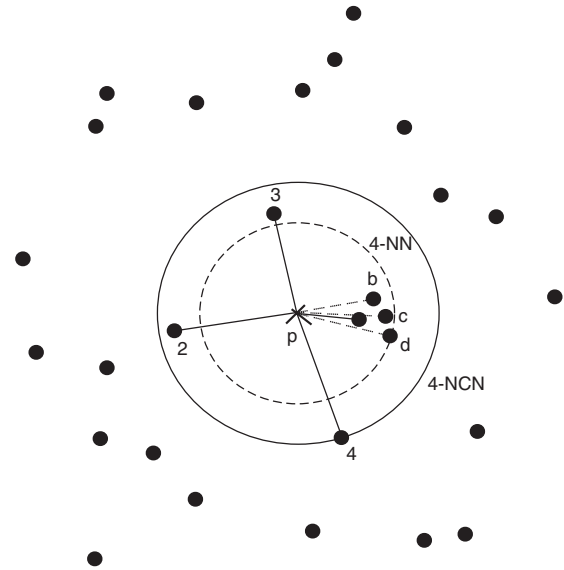


Fig. 1. Illustration of the NCN concept.

A set of prototypes is selected from the TS to guarantee their optimal geometrical distribution with respect to their NCNs. The use of the NCN of a given sample can provide local information about the shape of the probability class distribution, which depends on the nature and class of its NCNs, that is, on the nature of the prototypes in its surrounding area.

The rationale behind this approach is that the prototypes belonging to the same class are located in a neighbouring area and can be replaced by a single representative without significantly affecting the original boundaries. The main reason to employ the NCN instead of the NN is to benefit from its properties, i.e. that the NCN covers a bigger region than the NN, and that these neighbours are located in the area of influence around a given sample which is compensated in terms of their geometrical distribution.

Initially, all training objects are considered as prototypes. The algorithm attempts to replace a group of neighbouring prototypes of the same class by a representative. In order to decide which group of prototypes should be replaced, the NCN of each prototype  $p$  in the TS is computed until reaching a neighbour of a class different than the class of  $p$ . The prototype with the largest number of NCNs is defined as a representative of its corresponding group. Since this group lies in the area of influence defined by the NCN distribution, consequently, all its members can now be removed from the TS. Next, given the remaining prototypes, the algorithm updates the number of their neighbours (if some were previously eliminated) as belonging to the group of an already existing representative. This is repeated until there is no group of prototypes to be replaced by a representative. This basic scheme is called *MaxNCN*, and is presented in Algorithm 1.

**Algorithm 2** *Reconsistent*


---

```

for  $i = \text{each\_prototype}(TS)$  do
  neighbours_number[i] = 0
  neighbour = next_neighbour(i)
  while neighbour.class == i.class do
    neighbours_vector[i] = Id(neighbour)
    neighbours_number[i] ++
    neighbour = next_neighbour(i)
  end while
end for
while Max_neighbours() > 0 do
  EliminateNeighbours(id_Max_neighbours)
end while
count = 0
for  $i = \text{each\_prototype}(TS)$  do
  if Classify(i) != i.class then
    incorrect_class[count ++] = i
  end if
end for
for  $i = \text{each\_prototype}(\text{incorrect\_class}[])$  do
  neighbours_number_inc[i] = 0
  neighbour_inc = next_neighbour_inc(i)
  while neighbour_inc.class == i.class do
    neighbours_vector_inc[i] = Id(neighbour_inc)
    neighbours_number_inc[i] ++
    neighbour_inc = next_neighbour_inc(i)
  end while
end for
while Max_neighbours_inc() > 0 do
  EliminateNeighbours_inc(id_Max_neighbours_inc)
end while
AddCondensedIncToCondensedTS()

```

---

2.2. *Reconsistent*

The *Reconsistent* algorithm is an important modification of the *MaxNCN* algorithm towards obtaining a consistent condensed set [14]. The primary idea is that the consistency of a subset with respect to the TS should lead to a better classification. By using the *MaxNCN* algorithm some prototypes close to the decision boundaries are removed because of the order in which the instances are taken during the condensing process. The *Reconsistent* approach, proposed by us, tries to address this issue.

The procedure starts by applying the *MaxNCN* technique to a TS, yielding a reduced set, RS. Then, each object in the TS is tested by the 1-NN rule with respect to the current prototype set, RS. All misclassified objects form a new group, which is condensed using the RS set as the reference. In the end, this new condensed set is added to the RS, resulting in the final condensed set. The *Reconsistent* procedure is presented in Algorithm 2.

2.3. *LVQ*

In statistical pattern recognition, learning vector quantisation (LVQ) is a popular competitive learning algorithm, used in many applications. Its goal is to approximate the distribution of classes by using a reduced set of prototypes, while minimising the classification error [13]. It is, therefore, an adaptive condensing technique. In general, the classes can be described by a relatively small number of prototypes  $p_i$ , placed within each class region of the decision boundary by means of measures of neighbourhood.

In the initialisation step, the prototypes are placed in the TS, maintaining the same number of prototypes in each class. The class borders are represented in a piecewise linear way by segments of midplanes between the prototypes of neighbouring classes (the borders of the so-called Voronoi tessellations). This may seem to be a good strategy for approximating the class borders, using the fact that the average distances between the neighbouring prototypes should be the same on both sides of the borders.

Nevertheless, the optimal placement of the prototypes is not known. Therefore, their distances and their optimal cardinality,  $k$ , cannot be determined beforehand. To fix the value of  $k$  for each case, the number of prototypes has to be chosen first. A combination of editing and condensing is used for this purpose. Therefore, first the Wilson's editing ( $k$  of the  $k$ -NN used in the editing scheme, as explained in Section 4.1) is used, and then followed by the Hart's condensing algorithm [7]. The resulting size of the prototype set is close to ideal [4]. Using this number of prototypes, different values of  $k$  were tested, and the one leading to the highest classification accuracy was finally chosen. Different set sizes were used for each data set, in order to compare the results for different cases.

In the experiments, a variant of the original LVQ algorithm is used, namely the *optimised-learning-rate, OLVQ1*, [24]. The basis of this algorithm is the *LVQ1* [24] such that an individual learning rate  $\alpha_i(t)$  is selected for each prototype  $p_i$ . Several prototypes are assigned to each class such that the layout of the prototypes minimise approximately the misclassification errors in the 1-NN classification. The following equations define this process:

$$\begin{aligned}
 p_c(t+1) &= p_c(t) + \alpha_c(t)[x(t) - p_c(t)] && \text{if } x \text{ and } p_c \text{ are in the same class,} \\
 p_c(t+1) &= p_c(t) - \alpha_c(t)[x(t) - p_c(t)] && \text{if } x \text{ and } p_c \text{ are in different classes,} \\
 p_i(t+1) &= p_i(t) && \text{if } i \neq c,
 \end{aligned} \tag{1}$$

where  $x(t)$  is an input sample and  $p_c$  is the nearest  $p_i$  to  $x$ . In Ref. [24], the "optimal" values of  $\alpha_i(t)$  are determined by recursion as

$$\alpha_c(t) = \frac{\alpha_c}{1 + s(t)\alpha_c(t-1)}, \tag{2}$$



where  $s(t) = +1$  if the classification of the prototype  $p_c$  is correct and  $s(t) = -1$  if the classification is wrong. Therefore, it is necessary to stop the learning process after some “optimal” number of steps. Typically, the *OLVQI* may be stopped after 200 iterations.

#### 2.4. *MixtGauss*

*MixtGauss* is an adaptive condensing algorithm proposed by some of us in Ref. [15]. It is considered in the framework of mixture modelling by Gaussian distributions, while assuming a statistical independence of features. The prototypes are chosen as the mean vectors of the optimised Gaussians, whose mixtures are fit to model each of the classes. The details are given below.

In general, given a TS, we assume that each class follows a spatial distribution according to its class-conditional probability density function (*pdf*)  $P(\mathbf{x}|c_j)$  and the respective a priori probability  $P(c_j)$ ,  $c_j \in C$ . These class-conditional *pdfs* do not need to have a specific structure, in general. In practice, however, it is necessary to obtain their density estimations.

A natural way to deal with a density estimator is to consider a mixture density of modes. One approach to retain the capacity of  $P(\mathbf{x}|c_j)$  is by reflecting the local structure of the distribution by means of mixture modes  $P_m(\mathbf{x}|c_j)$ , where each mode is estimated by the product of probabilities in each feature as [26]:

$$\begin{aligned} P(\mathbf{x}|c_j) &= \sum_{m=1}^M \alpha_{m|j} P_m(\mathbf{x}|c_j) \\ &= \sum_{m=1}^M \alpha_{m|j} \prod_{k=1}^n N(x_k; \mu_{m|kj}, \sigma_{m|kj}), \end{aligned} \quad (3)$$

where  $M$  is the number of modes and  $\alpha_{m|j}$  is a priori probability of the  $m$ th mode in the class  $c_j$ .

In our condensing algorithm, each class in the TS is modelled by a probability distribution. A general shape of classes as well as the decision boundaries is maintained, when each class is described by a mixture of multivariate Gaussian distributions. By our additional assumption of the statistical independence of features, the  $m$ th component in a mixture modelling the class  $c_j$  is a multivariate Gaussian distribution expressed by a product of univariate normal distributions,  $N(\mu_{m|kj}, \sigma_{m|kj})$ . Note that this is equivalent to a multivariate elliptic Gaussian distribution  $N(\boldsymbol{\mu}_{m|j}, \text{diag}(\sigma_{m|j}))$ . The standard method used to fit finite mixture models to the observed data, hence to estimate the parameters of the class distributions, is the well-known EM algorithm [27,28]. This is a general maximum likelihood optimisation procedure for problems with hidden variables or missing data [29]. The final prototypes of the sought condensed set are the mean vectors of the Gaussian distributions determined by the EM algorithm.

Note that, although the EM algorithm is widely used, one needs to be aware of its drawbacks [30]. As a method working in local neighbourhoods, it is sensitive to initialisation, as the likelihood function of a mixture model is not unimodal. Another important problem in mixture modelling is the selection of the number of components. With too many components, the mixture may overfit the data, while a mixture with too few components may not be flexible enough to approximate the true underlying model.

---

#### Algorithm 3 *MixtGauss*

---

```
(* initialisation *)
for c = 1..number_of_classes do
  mean[c] = Calculate_Mean(c, TS)
  for g = 1..number_of_Gaussians_per_class do
    Gaussians[c, g] = mean[c] + Random_Disturbance()
  end for
end for
(* optimisation *)
repeat
  previous_Gaussians[c, g] = Gaussians[c, g]
  Gaussians[c, g] = EMstep()
  previous_accuracy = current_accuracy
  current_accuracy = Calculate_Accuracy()
  classes_improve = 0
  for c = 1..number_of_classes do
    if current_accuracy[c] > previous_accuracy[c] then
      classes_improve = classes_improve + 1
    else
      if current_accuracy[c] = previous_accuracy[c] then
        for g = 1..number_of_Gaussians_per_class do
          Gaussians[c, g] = previous_Gaussians[c, g]
        end for
      end if
    end if
  end for
until classes_improve == 0
```

---

##### 2.4.1. Algorithm details

The more Gaussian components are included in a mixture, the more accurate the representation of the classes and the decision boundaries is. Taking this into account,  $M$  Gaussians have to be specified to represent each class distribution. This number corresponds to the number of prototypes generated in each class for the final condensed set. Initially, each class is represented by  $M$  Gaussians located at the mean of that class. By adding random disturbances, different mean vectors are created and the Gaussians are shifted away. So, a mixture of Gaussians can be obtained per class. For each Gaussian, the initial variance is set up to  $\frac{1}{10}$  of the range in each dimension.

After the initialisation stage, the EM algorithm is used to determine an optimal location of the mixture of Gaussians. This iterative optimisation procedure converges to a

maximum likelihood estimate of the mixture parameters. Accordingly, to fit a mixture of Gaussians to each class, we iterate between the following two steps:

*E-step:* Compute the contributions of the prototypes  $\mathbf{x}^t$  in the class  $c_j$ , belonging to the set  $\{\mathbf{x}^1, \dots, \mathbf{x}^{N_j}\}$ , where  $N_j$  is the cardinality of the class  $c_j$ . The conditional pdf for the  $m$ th mode is

$$P_m(\mathbf{x}^t | c_j) = \frac{\alpha_{m|j} \prod_{k=1}^n N(x_k^t; \mu_{m|kj}, \sigma_{m|kj})}{\sum_{l=1}^M \alpha_{l|j} \prod_{k=1}^n N(x_k^t; \mu_{l|kj}, \sigma_{l|kj})} \quad (4)$$

and it is normalised such that  $\sum_{m=1}^M \alpha_{m|j} = 1$ . The multivariate Gaussians are represented as a product of univariate normal distributions, with the means  $\mu_{m|kj}$  and the standard deviations  $\sigma_{m|kj}$ .

*M-step:* Compute the parameters of the  $m$ th mode for each value  $\mathbf{x}^t$  that exists in the class  $c_j$ :

$$\alpha_{m|j} = \frac{1}{N_j} \sum_{t=1}^{N_j} P_m(\mathbf{x}^t | c_j),$$

$$\mu_{m|kj} = \frac{\sum_{t=1}^{N_j} P_m(\mathbf{x}^t | c_j) x_k^t}{\sum_{t=1}^{N_j} P_m(\mathbf{x}^t | c_j)}, \quad (5)$$

$$\sigma_{m|kj} = \frac{\sum_{t=1}^{N_j} P_m(\mathbf{x}^t | c_j) (x_k^t - \mu_{m|kj})^2}{\sum_{t=1}^{N_j} P_m(\mathbf{x}^t | c_j)}. \quad (6)$$

The EM iterative optimisation can cause an overlap between the Gaussian components from different classes, which will deteriorate the final classification accuracy. Therefore, this process should stop when no class yields an increase in performance with respect to the previous iteration. To address this problem, the consistency criterion is applied. After each iteration, the consistency criterion is estimated by calculating the 1-NN error over the TS by using the estimated means of the Gaussians as the current condensed set. The movement of the Gaussians is carried out while any class improves the classification rate with respect to the previous step. The complete technique is summarised in Algorithm 3.

### 3. Prototype-based classification methods to compare

The usefulness of condensed sets, optimised by the approaches discussed above, will be evaluated in a classification task. Traditionally, the 1-NN rule, assigning an unknown object to the class of its nearest neighbour in the condensed set, is used for this purpose. The 1-NN rule is employed, as it is often used in condensing schemes for the design of the condensed set. If the classes are represented as compact Gaussian-like clouds of similar spreads, then the 1-NN is expected to generalise well for a small condensed set. Otherwise, a large condensed set might be needed to represent the variability in the data. Alternatively, classifiers in the so-called dissimilarity spaces can be considered.

#### 3.1. Dissimilarity spaces

Remember that a TS  $X$  consist of  $N$  objects in a feature vector space. Let  $R = \{\mathbf{p}_1, \mathbf{p}_2, \dots, \mathbf{p}_r\}$  be a condensed set, called also representation set, of  $r$  optimised prototypes. Given a dissimilarity measure<sup>2</sup>  $d$ , one may consider a new representation based on the proximities to the set  $R$ . Every object  $\mathbf{x} \in X$  is then described by a vector of dissimilarities computed between  $\mathbf{x}$  and the prototypes from  $R$ , i.e. as  $D(\mathbf{x}, R) = [d(\mathbf{x}, \mathbf{p}_1), d(\mathbf{x}, \mathbf{p}_2), \dots, d(\mathbf{x}, \mathbf{p}_r)]$ . Hence, for a set  $X$ , it extends to an  $N \times r$  dissimilarity matrix  $D(X, R)$ .

The dissimilarity matrix  $D(X, R)$  is interpreted as a data-dependent mapping  $D(\cdot, R) : X \rightarrow \mathbb{R}^r$  from the representation  $X$  to a *dissimilarity space*, defined by the set  $R$ . This is a vector space, in which each dimension corresponds to a dissimilarity to a prototype from  $R$ ,  $D(\cdot, \mathbf{p}_i)$ . The advantage of this representation is that any traditional classifier operating in vector spaces can now be used [20,16,17].

A vector  $D(\cdot, \mathbf{p}_i)$  of dissimilarities to the prototype  $\mathbf{p}_i$  can be interpreted as a feature. This follows since the dissimilarities to  $\mathbf{p}_i$  will differ depending on the class membership. If the measure is metric and the dissimilarity  $d(\mathbf{p}_i, \mathbf{p}_j)$  is small, then  $d(\mathbf{x}, \mathbf{p}_i) \approx d(\mathbf{x}, \mathbf{p}_j)$  for other object  $\mathbf{x}$ , as it is guaranteed by the backward triangle inequality [17]. As a result, only one of them can be chosen as a prototype. However, a small number of prototypes might be insufficient to represent the data variability if the classes have different spreads (in terms of the dissimilarities). This occurs when, in a feature space one class is represented by a compact cloud of points, while the other class is elongated. It can be, therefore, more beneficial to inspect the vectors of dissimilarities to the entire set  $R$ , instead of looking at the nearest neighbour only. If the objects  $\mathbf{x}$  and  $\mathbf{y}$  are similar, then their dissimilarity vectors  $D(\mathbf{x}, R)$  and  $D(\mathbf{y}, R)$  are expected to be correlated. Hence they should lie close in a dissimilarity space (as measured e.g. by an Euclidean distance in this space). A linear or quadratic classifier in this dissimilarity space might be of interest.

#### 3.2. Normal-density based classifiers in dissimilarity spaces

In the case of small condensed sets or non-representative TS, a better generalisation can be achieved by a classifier built in a dissimilarity space than by the 1-NN rule. Many traditional classifiers can be applied there [19–21,17]. Linear and quadratic functions are weighted linear (quadratic) combinations of the dissimilarities  $d(\mathbf{x}, \mathbf{p}_i)$  between a given object  $\mathbf{x}$  and the prototypes  $\mathbf{p}_i$ . Although the classifiers are trained on  $D(\cdot, R)$ , the weights are still optimised on the complete TS.

<sup>2</sup> In general, a non-negative dissimilarity measure should express a degree of commonality between pairs of objects. It should be zero for two identical objects, take small values for similar objects and large values for objects that differ. It does not need to be metric.

Bayesian classifiers, i.e. linear and quadratic normal density based classifiers, tend to perform well in dissimilarity spaces [19,18,20,21,31]. This is especially true, for a summation-based dissimilarity measure, summing over a number of components with similar variances. The reason is that such dissimilarities will be approximately normally distributed thanks to the central limit theorem (if one or few variances are dominant, then they will approximate the  $\chi^2$  distribution) [17]. In our case, the vectorial data will be normalised to unit variances, hence the resulting Euclidean distances tend to be approximately normally distributed. One may, therefore, expect a good performance of such decision rules, if the problem is linearly or quadratically separable. Since in our case the Euclidean distance is a non-linear transformation normalised features, a linear classifier in a dissimilarity space is a non-linear classifier in the underlying (original) space [17], similarly to the reasoning for support vector machines [32]. So, if the Euclidean distance is informative for the classification problem in a feature vector space, the local and non-linear 1-NN rule in this space may be replaced by a linear classifier in the Euclidean-distance dissimilarity space.

For a two-class problem, a linear normal density based classifier (NLC), defined by the condensed set  $R$ , is given by

$$f(D(x, R)) = [D(x, R) - \frac{1}{2}(\mathbf{m}_{(1)} + \mathbf{m}_{(2)})]^T \times S^{-1}(\mathbf{m}_{(1)} - \mathbf{m}_{(2)}) + \log \frac{p_{(1)}}{p_{(2)}}, \quad (7)$$

and the quadratic function (NQC) becomes

$$f(D(x, R)) = \sum_{i=1}^2 \frac{(-1)^i}{2} (D(x, R) - \mathbf{m}_{(i)})^T \times S_{(i)}^{-1} (D(x, R) - \mathbf{m}_{(i)}) + \log \frac{p_{(1)}}{p_{(2)}} + \frac{1}{2} \log \frac{|S_{(1)}|}{|S_{(2)}|}, \quad (8)$$

where  $\mathbf{m}_{(1)}$  and  $\mathbf{m}_{(2)}$  are the mean vectors,  $S$  is the sample covariance matrix (average of the class covariance matrices) and  $S_{(1)}$  and  $S_{(2)}$  are the estimated class covariance matrices, all computed in the dissimilarity space  $D(X, R)$ .  $p_{(1)}$  and  $p_{(2)}$  are the class prior probabilities. When the covariance matrices become singular, they can be regularised. In the implementation of these classifiers for multi-class problems, the normal-density functions are estimated per each class and the final decision is based on the maximum a posteriori probability [33].

Given equal class priors, the NLC is equivalent to the Fisher linear discriminant (FLD) obtained by maximising the *Fisher criterion*, i.e.  $\max_{\mathbf{w}} (\mathbf{w}^T S_B \mathbf{w}) / (\mathbf{w}^T S_W \mathbf{w})$ , where  $S_B$  is the between-class scatter and  $S_W = S$  is the within-class scatter; see Ref. [34] for details. For a representation  $D(T, R)$ , the FLD is found as

$$f(D(x, R)) = (\mathbf{m}_1 - \mathbf{m}_2)^T S_W^{-1} D(x, R) - \frac{1}{2} (\mathbf{m}_1 + \mathbf{m}_2)^T S_W^{-1} (\mathbf{m}_1 - \mathbf{m}_2). \quad (9)$$

Table 1  
Data sets used in the experiments

Data set	No. classes	No. features	TS size	Test set size
Cancer	2	9	546	137
Diabetes	2	8	614	154
Glass	6	9	171	43
Heart	2	13	216	54
Liver	2	6	276	69
Vehicle	4	18	677	169
Vowel	11	10	422	106
Wine	3	13	142	36
Phoneme	2	5	4323	1080
Satimage	6	36	5148	1287
Texture	11	40	4400	1100

A multi-class FLD is derived in the one-against-all strategy [33].

#### 4. Experimental results and discussion

We will compare some methods of determining a condensed set  $R$  in combination with three prototype-based classification strategies. These is done on Euclidean distance representations  $D(X, R)$  derived in feature spaces. The four condensing algorithms are the ones described before: *MaxNCN*, *Reconsistent*, *LVQ* and *MixtGauss*. The first classification method is the 1-NN rule based on the set  $R$ . The other two methods are the FLD and NQC trained in dissimilarity spaces  $D(X, R)$ .

##### 4.1. Data description and experimental setup

Eleven real data sets (see Table 1) are taken from the UCI Repository [35] to assess the behaviour of the algorithms presented here. All data sets are first normalised by a unit variance, and then the Euclidean distances are computed. The results are compared in terms of both accuracy and reduction of the TS size.

The *MaxNCN*, *Reconsistent* and *MixtGauss* algorithms need to be applied to overlap-free (no overlap between different class regions) data sets. Therefore, as a general rule, and according to the previously published results [10,5], the Wilson's editing is considered to properly remove a possible overlap between the classes. In our experiments, the parameter  $k$  involved (of the  $k$ -NN rule) is obtained by performing a five-fold cross-validation experiment using only the TS and computing the average classification accuracies for different values of  $k$ . The best edited set (including the non-edited TS) is thus selected as the input for the different condensing schemes.

##### 4.2. Quantitative results

Tables 2–4 report the classification accuracy of the 1-NN, the FLD and the NQC, respectively, obtained by

Table 2  
Classification accuracy of the 1-NN (in %)

	Original	MaxNCN	Reconsistent	LVQ	MixtGauss
Cancer	95.2	92.0 [95.2 ≡ 26]	93.4 [93.8 ≡ 34]	96.9 [99.3 ≡ 4]	96.5 [99.3 ≡ 4]
Diabetes	70.1	67.7 [87.0 ≡ 80]	71.8 [81.8 ≡ 112]	76.6 [98.1 ≡ 12]	75.0 [91.9 ≡ 50]
Glass	69.7	64.1 [59.1 ≡ 70]	69.4 [46.8 ≡ 91]	72.6 [71.9 ≡ 48]	64.5 [75.4 ≡ 42]
Heart	76.3	69.3 [83.8 ≡ 35]	74.8 [77.8 ≡ 48]	84.1 [89.8 ≡ 22]	81.9 [93.5 ≡ 14]
Liver	62.6	58.9 [76.8 ≡ 64]	61.0 [68.5 ≡ 87]	65.6 [93.5 ≡ 18]	63.7 [97.1 ≡ 8]
Vehicle	69.4	64.4 [60.3 ≡ 269]	66.0 [49.2 ≡ 344]	72.2 [61.6 ≡ 260]	69.6 [70.5 ≡ 200]
Vowel	98.1	90.8 [74.4 ≡ 108]	95.5 [71.8 ≡ 119]	80.8 [76.5 ≡ 99]	90.7 [76.5 ≡ 99]
Wine	94.4	91.6 [93.0 ≡ 10]	92.1 [90.9 ≡ 13]	95.5 [95.8 ≡ 6]	96.0 [97.9 ≡ 3]
Phoneme	71.0	69.1 [91.6 ≡ 365]	70.2 [88.9 ≡ 480]	77.3 [93.1 ≡ 300]	73.3 [93.1 ≡ 300]
Satimage	83.0	79.2 [86.5 ≡ 693]	80.3 [82.7 ≡ 890]	82.2 [88.3 ≡ 600]	81.3 [94.8 ≡ 270]
Texture	98.9	95.3 [89.3 ≡ 470]	97.3 [87.8 ≡ 535]	85.9 [99.5 ≡ 22]	97.2 [90.0 ≡ 440]
Average	80.8	76.6 [81.5]	79.2 [76.4]	80.9 [87.9]	80.9 [89.1]

The corresponding reduction rates of the TS size (in %) and the equivalent number of prototypes are given in brackets.

Table 3  
Classification accuracy of the FLD in dissimilarity spaces (in %)

	Original	MaxNCN	Reconsistent	LVQ	MixtGauss
1 Cancer	85.50	96.48 [95.2]	96.48 [93.8]	97.22 [98.2]	96.92 [96.3]
2 Diabetes	74.23	76.05 [87.0]	77.09 [81.8]	77.08 [98.7]	76.70 [87.0]
3 Glass	53.88	66.53 [59.1]	67.49 [46.8]	77.17 [61.4]	69.77 [93.0]
4 Heart	60.75	80.72 [83.8]	78.49 [77.8]	84.79 [98.1]	82.60 [98.2]
5 Liver	68.35	68.96 [76.8]	67.83 [68.5]	70.70 [92.8]	68.67 [93.5]
6 Vehicle	77.19	79.31 [60.3]	79.79 [49.2]	80.96 [67.5]	80.02 [76.4]
7 Vowel	41.03	95.76 [74.4]	96.67 [71.8]	91.42 [76.5]	94.28 [76.5]
8 Wine	33.15	98.28 [93.0]	98.30 [90.9]	98.87 [95.8]	98.90 [95.8]
9 Phoneme	70.62	69.75 [91.6]	70.52 [88.9]	77.11 [98.2]	73.94 [99.8]
10 Satimage	****	24.49 [86.5]	23.90 [82.7]	82.11 [95.3]	82.20 [95.9]
11 Texture	****	37.33 [89.3]	27.56 [87.8]	99.00 [92.5]	99.13 [93.8]
Average (1–11)	****	72.2 [81.5]	71.3 [76.4]	85.1 [88.6]	83.9 [91.5]
Average (1–9)	62.7	81.3	81.4	83.9	82.4

The corresponding reduction rates of the TS size (in %) are given in brackets. Averages over all data sets (1–11) and data sets excluding Satimage and Texture (1–9), are also given.

Table 4  
Classification accuracy of the NQC in dissimilarity spaces (in %)

	Original	MaxNCN	Reconsistent	LVQ	MixtGauss
1 Cancer	80.23	96.19 [95.2]	96.48 [93.8]	97.22 [96.3]	96.34 [96.3]
2 Diabetes	65.62	75.39 [87.0]	74.74 [81.8]	77.09 [97.1]	76.70 [98.7]
3 Glass	40.15	53.81 [59.1]	54.40 [46.8]	65.10 [89.5]	64.11 [82.5]
4 Heart	55.56	82.97 [83.8]	80.00 [77.8]	82.95 [97.2]	82.21 [88.9]
5 Liver	46.67	62.63 [76.8]	60.00 [68.5]	67.79 [85.5]	65.20 [89.1]
6 Vehicle	27.78	60.27 [60.3]	54.37 [49.2]	74.81 [76.4]	74.58 [91.1]
7 Vowel	10.20	31.96 [74.4]	29.41 [71.8]	94.59 [92.2]	94.10 [92.2]
8 Wine	33.15	97.73 [93.0]	96.05 [90.9]	97.76 [97.9]	99.43 [95.8]
9 Phoneme	73.26	73.08 [91.6]	73.22 [88.9]	76.11 [95.4]	73.33 [95.4]
10 Satimage	****	79.40 [86.5]	79.17 [82.7]	81.93 [98.3]	81.62 [97.1]
11 Texture	****	96.31 [89.3]	96.51 [87.8]	97.58 [95.0]	97.53 [96.3]
Average (1–11)	****	73.6 [81.5]	72.2 [76.4]	83.0 [92.8]	82.3 [93.0]
Average (1–9)	48.1	70.5	68.7	81.5	80.7

The corresponding reduction rates of the TS size (in %) are given in brackets. Averages over all data sets (1–11) and data sets excluding Satimage and Texture (1–9), are also given.

using different condensed sets. For each condensing method, the average performance values, computed over the 11 data sets are also included. There are no results of the FLD and

the NQC for the *Satimage* and *Texture* sets (for the original TS), because of their high memory requirements during training [33]. The average performances excluding these



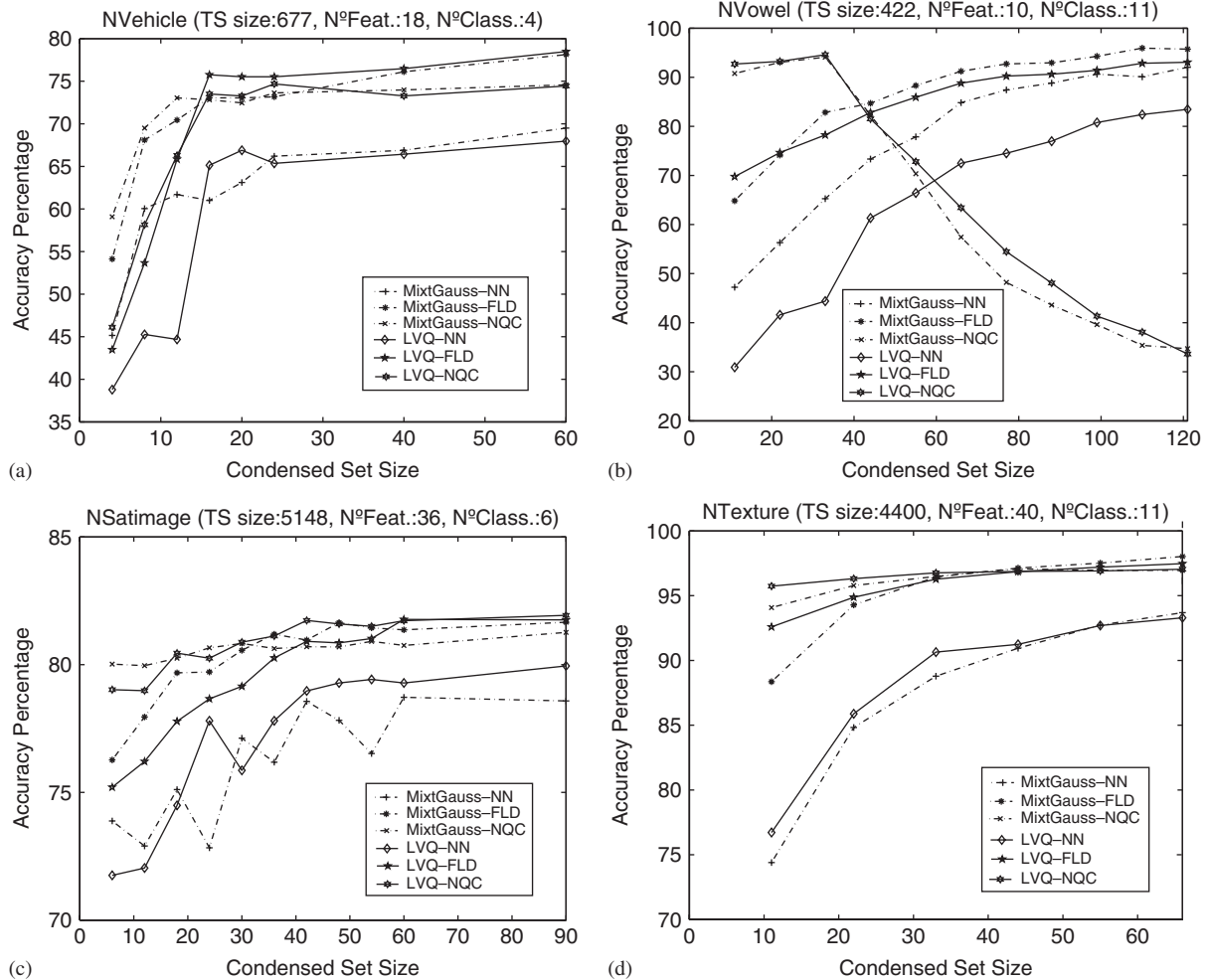


Fig. 2. Trade-off between the resulting condensed set sizes and classification accuracy for different algorithms. The following data sets are shown: *Vehicle* (a); *Vowel* (b); *Satimage* (c) and *Texture* (d).

data sets are also presented. Since the adaptive techniques, the *MixtGauss* and the *LVQ*, allow one to determine a condensed set of any size, only a few results are shown in the tables. These correspond to the best accuracy values, found in the sets between one prototype per class and the size determined by the *MaxNCN* algorithm (as the *MaxNCN* condensed set is always smaller than the *Reconsistent* condensed set).

The first observation is that the *LVQ* and the *MixtGauss* increase the average accuracy in comparison to the original TS accuracy. The increase is larger by using the *FLD* and the *NQC* than by using the *1-NN* rule (where the increase is quite small). Anyway, it is significant that in all these cases, the performance is increased despite the high reduction of the TS (see Tables 2–4). Remember that although the *FLD* and the *NQC* are based on the distances to the prototypes, they still make use of the entire TS in training. This is the reason why the classification accuracy increases so much when using any condensed set, in comparison to the results by using the original TS.

By using the *FLD*, the classification accuracy obtained for the *MaxNCN* and *Reconsistent* procedures is very similar. The *MixtGauss* leads to a higher accuracy than both of them. The best performance of the *FLD* relies on the *LVQ*. A similar pattern can be observed for the *NQC*. This differs, however, for the *1-NN* rule, especially in the case of the *Reconsistent* algorithm. The smallest classification accuracy is reached for the *MaxNCN*. The *Reconsistent* gives better results, and the best values are obtained for the *LVQ* and the *MixtGauss*.

The training reduction rates, being the percentage by which the original training size is reduced, are provided in square brackets in Tables 2–4. Since the *MaxNCN* and the *Reconsistent* techniques determine the optimal number of prototypes directly, their reduction rates are identical for the three classification methods used. Additionally, for each condensing method, the average reduction values computed over the 11 data sets are also included.

In relation to the training reduction rates, the average is higher for the *MaxNCN* than for the *Reconsistent*. This is

to be expected as in the second algorithm, some instances are added to the MaxNCN condensed set to improve the classification accuracy. It can also be observed that the average TS reduction rates are higher for the adaptive *LVQ* and *MixtGauss* techniques than for the selective *MaxNCN* and *Reconsistent* procedures. The highest reduction rate is the one of *MixtGauss*. This holds for the three prototype-based classification strategies used here.

In relation to the classification methods considered, in general, we can observe that the *FLD* and the *NQC* perform better than the *1-NN* rule, when they rely on the condensed sets determined by the adaptive schemes. Comparing both algorithms, the one with the best performance corresponds to the *FLD*. On the other hand, the *MaxNCN* and the *Reconsistent* procedures have their best classification accuracy with the *1-NN* rule. As designed, the *Reconsistent* leads to a higher nearest neighbour performance than the *MaxNCN*.

These results can be better analysed by studying the graphs for each data set (see an example for four data sets in Fig. 2). The best trade-off results (between the accuracy and the reduced TS size) are the ones for the *FLD*. The results for the *NQC* are similar, and the ones for the *1-NN* rule are significantly worse. In relation to the condensing algorithms, the *LVQ* and the *MixtGauss* lead to very similar results. In general, the points corresponding to the *MaxNCN-FLD* and the *Reconsistent-FLD* results are near the functions representing the *MixtGauss-FLD* and the *LVQ-FLD* results, similarly as the points representing the *MaxNCN-NQC* and the *Reconsistent-NQC* results are close to the points describing the *MixtGauss-NQC* and the *LVQ-NQC* results. This observation does not hold for the *MaxNCN-NN* and the *Reconsistent-NN*, neither the *MixtGauss-NN* and the *LVQ-NN*, as there exists a big difference in favour of the *MixtGauss* and the *LVQ* for some data sets.

In Fig. 2, the behaviour of the algorithms investigated here is shown for four data sets. Only the results for the smallest reduced sets for each data are plotted, in order to distinguish the plots. The first graph, Fig. 2(a), refers to the *Vehicle* data set. It can be observed there that the best classification accuracy is reached by the *MixtGauss-NQC* for the smallest condensed sets. For sets with more than three instances per class, the best accuracy is reached by the *LVQ-FLD*. It is to note that the performance of the *MixtGauss-FLD* is always close to the best case.

In Fig. 2(b) the results for the *Vowel* data set are shown. The three smallest condensed sets represented (1, 2 and 3 instances per class) lead to the best performance for the combination of the *LVQ-NQC* and the *MixtGauss-NQC* algorithms. However, for larger sizes, the classification accuracy suddenly decreases. It is due to a large number of classes in relation to the number of prototypes, which effectively translates into the dimension of the dissimilarity space. As a result, the number of instances per class is too small, and as the size of the condensed set increases, class covariance matrices are badly estimated. Hence, the classification accuracy is diminished. The next best results belong

to the *MixtGauss-FLD*, the *LVQ-FLD* and the *MixtGauss-NN*, in that order.

In Fig. 2(c), the results for the *Satimage* data set are shown. The most reduced sets lead to a better performance by using the *MixtGauss-NQC*, the *LVQ-NQC*, the *MixtGauss-FLD* and the *LVQ-FLD* methods. The results for the data set *Texture* are shown in Fig. 2(d). The best ones are obtained by the following combinations: the *LVQ-NQC*, the *MixtGauss-NQC*, the *LVQ-FLD* and the *MixtGauss-FLD*. As it is not plotted in this paper, in order to save space, it could be important to know the next point. The larger graphs for *Satimage* and *Texture* data sets show that the classification accuracy also decreases around the size of 300 instances for the *MixtGauss-FLD* and the *LVQ-FLD*. It is due to similar reasons as mentioned for the *Vowel* data.

## 5. Conclusions

This paper compares several prototype reduction techniques, including both selecting and generating schemes, in connection with performance of classifiers defined on the corresponding sets of prototypes. We use the following classifiers: the *1-NN* rule, and linear and quadratic prototype-based classifiers defined in Euclidean-distance dissimilarity spaces. The purpose of this experimental study is to discuss a number of prototype optimisation methods with respect to their ability of maintaining small prototype sets and high classification performance.

From our investigations it can be concluded that there is no significant difference between the *LVQ* and the *MixtGauss* prototype generation techniques. However, when comparing them to the *MaxNCN* and *Reconsistent* methods, the first two algorithms tend to give both larger reduction of the TS (i.e. smaller condensed sets) and lead to higher classification accuracy. This means that the *LVQ* and the *MixtGauss* are good as condensing algorithms since they yield a good performance and allow one to control the number of prototypes. Additionally, they represent a good trade-off between the classification accuracy and the reduction rates.

In relation to the prototype-based decision rules, in general, the classification accuracy obtained by the *FLD* tends to be higher than that of the *NQC*, both higher than the accuracy of the *1-NN* rule, independently of the considered condensing algorithms. This is striking as the prototype sets are often optimised to guarantee that the *1-NN* rule performs well. Our explanation is that the *FLD* and the *NQC*, although based on the condensed sets, still make use of the entire TS for determining the decision boundary. The *1-NN* rule, when applied as a nearest prototype rule, discards other training objects.

What is important, is the fact that for the evaluation of new objects the computational complexity of the *FLD* is  $O(Jr)$ . This is similar to the computational complexity of the *1-NN* rule,  $O(r)$ , provided that  $J$ , the number of classes, is small. For a large number of classes, the computational cost of the *FLD* increases. The *FLD* may, therefore, be an

alternative to the 1-NN rule for the prototype-based classification in normalised feature vector spaces. The *NQC* is much more computationally heavy, so it can only be advantageous for small condensed sets and a small number of classes. In fact, the *NQC* was used here, to show that a more complex classifier than a linear one is not necessarily yielding a better performance in dissimilarity spaces resulting from normalised vectorial data. The effect of feature normalisation plays a role as otherwise the *NQC* might have been preferred to the *FLD*. In general, the good performance of the *FLD* in dissimilarity spaces can be explained by the use of Euclidean distance relying on *standardised* features (hence approximate normal distributions for the resulting distances) and non-complicated class structures for which this distance measure is discriminative. The later means that the classes do not differ much in the range of their average within-class distances.

Our final conclusion is that the adaptive techniques, the *LVQ* and the *MixtGauss*, combined with the *FLD*, offer the best trade-off between the reduction rate, computational cost and the classification performance in normalised vector spaces for the considered problems.

In perspective, our study opens a new possibility for using linear (or more complex) classifiers in dissimilarity spaces derived from normalised feature vector spaces instead of the *1-NN* rule, both defined by the same optimised prototype sets. Since a linear function in such a dissimilarity space is a non-linear function in the original vector space, the distance measure or its non-linear (sigmoidal, logarithmic) transformation can be a way to incorporate the non-linearity aspects of the data. To our knowledge, these aspects have not been considered among the researchers studying the condensing techniques, yet. Although the *FLD* was chosen here, other linear functions can be studied, such as a logistic classifier or a hyperplane defined by a linear, sparse or non-sparse, programming procedure. Additionally, a linear (or polynomial) support vector machine (SVM) [32] can be trained in dissimilarity spaces. Note, however, that in such a case the SVM will rely on support vectors which are objects found in a complete dissimilarity space, hence they rely on the distances to *all* original objects. So, the SVM will not work as a prototype selection. If however, good prototypes are found by other techniques, the linear SVM will determine the optimal large-margin hyperplane in the corresponding dissimilarity space. Future research may focus on comparing a (non-)linear SVM built in the original feature vector space and the SVM in a dissimilarity space defined by a small set of optimised prototypes.

## Acknowledgments

This work has been partially supported by grants DPI2001-2956-C02-02 and TIC2003-08496 from the Spanish CICYT and project IST-2001-37306 from the European Union, as well as by the Dutch Organisation for Scientific Research (NWO).

## References

- [1] D.L. Wilson, Asymptotic properties of nearest neighbor rules using edited data sets, *IEEE Trans. Syst. Man Cybern.* 2 (1972) 408–421.
- [2] D.W. Aha, D. Kibler, M.K. Albert, Instance-based learning algorithms, *Mach. Learn.* 6 (1) (1991) 37–66.
- [3] B.V. Dasarthy, Nearest Neighbor (NN) Norms: NN Pattern Classification Techniques, IEEE Computer Society Press, Los Alamitos, CA, 1990.
- [4] P.A. Devijver, J. Kittler, *Pattern Recognition: A Statistical Approach*, Prentice-Hall, Englewood Cliffs, NJ, 1982.
- [5] D.R. Wilson, T.R. Martinez, Reduction techniques for instance-based learning algorithms, *Mach. Learn.* 38 (3) (2000) 257–286.
- [6] B.V. Dasarthy, Minimal consistent subset (MCS) identification for optimal nearest neighbor decision systems design, *IEEE Trans. Syst. Man Cybern.* 24 (1994) 511–517.
- [7] P.E. Hart, The condensed nearest neighbor rule, *IEEE Trans. Inform. Theory* 14 (5) (1968) 505–516.
- [8] I. Tomek, Two modifications of CNN, *IEEE Trans. Syst. Man Cybern.* 6 (1976) 769–772.
- [9] G.T. Toussaint, B.K. Bhattacharya, R.S. Poulsen, The application of voronoi diagrams to nonparametric decision rules, in: L. Billard (Ed.), *Computer Science and Statistics: The Interface*, Elsevier Science, Amsterdam, 1985.
- [10] M.C. Ainslie, J.S. Sánchez, Space partitioning for instance reduction in lazy learning algorithms, *Second Workshop on Integration and Collaboration Aspects of Data Mining, Decision Support and Meta-Learning*, 2002, pp. 13–18.
- [11] C.L. Chang, Finding prototypes for nearest neighbor classifiers, *IEEE Trans. Comput.* 23 (1974) 1179–1184.
- [12] C.H. Chen, A. Jóźwik, A sample set condensation algorithm for the class sensitive artificial neural network, *Pattern Recognition Lett.* 17 (1996) 819–823.
- [13] T. Kohonen, *Self-Organizing Maps*, Springer, Berlin, 1995.
- [14] M. Lozano, J.S., Sánchez, F. Pla, Using the geometrical distribution of prototypes for training set condensing, in: R. Conejo, M. Urretavizcaya, J.L. Pérez de la Cruz (Eds.), *Current Topics in Artificial Intelligence, Lecture Notes in Artificial Intelligence*, vol. 3040, Springer, Berlin, 2004, 618–627.
- [15] M. Lozano, J.M. Sotoca, J.S. Sánchez, F. Pla, An adaptive condensing algorithm based on mixtures of Gaussians, in: J. Vitrià et al. (Eds.), *Setè Congrés Català d'Intel·ligència Artificial, Barcelona (Spain), Frontiers in Artificial Intelligence and Applications*, vol. 113, IOS Press, Amsterdam, 2004, pp. 225–232.
- [16] E. Pełkalska, R.P.W. Duin, P. Paclík, Prototype selection for dissimilarity-based classifiers, *Pattern Recognition* 39 (2) (2006) 189–208.
- [17] E. Pełkalska, R.P.W. Duin, *The Dissimilarity Representation for Pattern Recognition. Foundations and Applications*, World Scientific, Singapore, 2005.
- [18] P. Paclík, R.P.W. Duin, *Classifying spectral data using relational representation*, *Spectral Imaging Workshop*, Graz, Austria, 2003.
- [19] P. Paclík, R.P.W. Duin, Dissimilarity-based classification of spectra: computational issues, *Real Time Imaging* 9 (4) (2003) 237–244.
- [20] E. Pełkalska, R.P.W. Duin, Dissimilarity representations allow for building good classifiers, *Pattern Recognition Lett.* 23 (8) (2002) 943–956.
- [21] E. Pełkalska, P. Paclík, R.P.W. Duin, A generalized kernel approach to dissimilarity based classification, *J. Mach. Learn. Res.* 2 (2) (2002) 175–211.
- [22] T.M. Cover, P.E. Hart, Nearest neighbor pattern classification, *IEEE Trans. Inform. Theory* 13 (1) (1967) 21–27.
- [23] J.S. Sánchez, F. Pla, F.J. Ferri, On the use of neighbourhood-based non-parametric classifiers, *Pattern Recognition Lett.* 18 (1997) 1179–1186.
- [24] T. Kohonen, J. Hynninen, J. Kangas, J. Laarksonen, K. Torkkola, *LVQ\_PAK: the learning vector quantization program package*,

- Helsinki University of Technology, 1996 (<http://www.cis.hut.fi/research/som-research/nnrc-programs.shtml>).
- [25] B.B. Chaudhuri, A new definition of neighbourhood of a point in multi-dimensional space, *Pattern Recognition Lett.* 17 (1996) 11–17.
- [26] J. Novovičová, P. Pudil, J. Kittler, Divergence based feature selection for multimodal class densities, *IEEE Trans. Pattern Anal. Mach. Intell.* 18 (1996) 218–223.
- [27] G. McLachlan, K. Basford, *Mixture Models: Inference and Application to Clustering*, Marcel Dekker, New York, 1988.
- [28] G. McLachlan, T. Krishnan, *The EM Algorithm and Extensions*, Wiley, New York, 1997.
- [29] A.P. Dempster, N.M. Laird, D.B. Rubin, Maximum likelihood from incomplete data via the EM algorithm, *J. R. Statist. Soc. B* 39 (1977) 1–38.
- [30] M.A.T. Figueiredo, A.K. Jain, Unsupervised learning of finite mixture models, *IEEE Trans. Pattern Anal. Mach. Intell.* 24 (2002) 381–396.
- [31] E. Pełalska, R.P.W. Duin, S. Günter, H. Bunke, On not making dissimilarities euclidean, *Structural Syntactic and Statistical Pattern Recognition*, Lecture Notes in Computer Science, vol. 3138, Springer, Berlin, 2004, pp. 1145–1154.
- [32] N. Cristianini, J. Shawe-Taylor, *An Introduction to Support Vector Machines*, Cambridge University Press, Cambridge, 2000.
- [33] R.P.W. Duin, P. Juszczak, D. de Ridder, P. Paclík, E. Pełalska, D.M.J. Tax, PR-tools, a Matlab toolbox for pattern recognition, 2004 (<http://www.prtools.org>).
- [34] R.O. Duda, P.E. Hart, D.G. Stork, *Pattern Classification*, second ed., Wiley, New York, 2001.
- [35] C.J. Merz, P.M. Murphy, UCI Repository of Machine Learning DB, Department of Information and Computer Science, University of California, 1998 (<http://www.ics.uci.edu/~mlearn>).

**About the Author**—M. LOZANO is a Lecturer at the Computer Science and Systems Engineering Department at the University of Zaragoza (Spain). She received a degree in Computing Engineering from the University Jaume I (Castellón, Spain) in 2000. Her Ph.D. dissertation is focused on the selection and generation of prototypes in learning processes.

**About the Author**—J.M. SOTOCA is an Assistant Lecturer at the Department of Programming Languages and Computer Systems of the Jaume I University. He received a degree and a Ph.D. in Physics in 1996 and 2001, respectively. His main research interests lie in the area of pattern recognition and computer vision. He has published more than 30 scientific papers in conferences, books and journals.

**About the Author**—J.S. SÁNCHEZ is an Associate Professor in the Department of Programming Languages and Information Systems at Universitat Jaume I (Castelló de la Plana, Spain) since 1992, and is currently the head of the Pattern Recognition Section in the Computer Vision Lab. He received a B.Sc. in Computer Science from the Universidad Politécnic de Valencia in 1990 and a Ph.D. in Computer Science Engineering from Universitat Jaume I in 1998. He is a member of IEEE, IAPR, and ECCAI. He is author or co-author of more than 80 scientific publications, co-editor of two books and guest editor of three special issues in international journals. His current research interests lie in the areas of pattern recognition and machine learning, including non-parametric classification, data reduction, data complexity analysis, and ensembles of classifiers.

**About the Author**—F. PLA is a Full Professor at the Department of Programming Languages and Information Systems of University Jaume I (Castellón, Spain), and is currently the group coordinator in the Computer Vision Lab. He received a degree and a Ph.D. in Physics from the University of Valencia in 1989 and 1993, respectively. He has authored more than 90 scientific papers in the fields of Computer Vision and Pattern Recognition. He has also been co-editor of two books and acted as reviewer for several International Journals in the field of Computer Vision and Pattern Recognition. He has been visiting scientist at Silsoe Research Institute, the University of Surrey and the University of Bristol in UK, CEMAGREF in France, the University of Genoa in Italy, and the Instituto Superior Técnico of Lisbon in Portugal. His current research interests are colour and spectral image analysis, visual motion analysis, active vision and pattern recognition techniques applied to image processing. He is member of the International Association for Pattern Recognition.

**About the Author**—E. PEŁALSKA received an M.Sc. degree in computer science from the University of Wrocław, Poland. In 1998, she became a research fellow working on an applied project at Delft University of Technology, The Netherlands, where in 1999 she started her research on dissimilarity-based learning methods. She obtained a cum-laude Ph.D. degree in January 2005. Currently, she is a postdoc at the same university and a honorary visiting research associate at the University of Manchester, UK.

**About the Author**—R.P.W. DUIN studied applied physics at Delft University of Technology in the Netherlands. In 1978 he received the Ph.D. degree for a thesis on the accuracy of statistical classifiers. In his research he focuses on the foundation of pattern recognition systems aiming at representations suitable for automatic learning.