# Classes and clusters in data analysis

A.M. Rubinov, N.V. Sukhorukova and J. Ugon<sup>\*</sup>

CIAO and School of Information Technology and Mathematical Sciences University of Ballarat P.O. Box 663 Ballarat, Victoria 3353 Australia

## Abstract

We discuss the relation between classes and clusters in datasets with given classes. We examine the distribution of classes within obtained clusters, using different clustering methods which are based on different techniques. We also study the structure of the obtained clusters. One of the main conclusions, obtained in this research is that the notion purity cannot be always used for evaluation of accuracy of clustering techniques.

*Key words:* data mining, dataset with classes, point-based clustering, optimization clustering model, comparison of classes and clusters.

# 1 Introduction

The goal of this paper is to study relations between classes and clusters in real-world datasets. This research topic is important in the area of data classification. There are two types of data classification: supervised and unsupervised.

In supervised classification the algorithms are provided with the data points (observations) and the labels to indicate the class each observation belongs to. The task of supervised classification is to develop decision rules for prediction of the class for any new observation from its attribute values. In other words the algorithms should predict the label for each new observation.

In unsupervised classification (clustering) the algorithms work just with data points, there is no division into classes. The task of clustering algo-

Preprint submitted to Elsevier Science

<sup>\*</sup> Corresponding author: jugon@ballarat.edu.au

rithms is to divide data items into several homogeneous groups (clusters): similar observations are grouped in the same cluster. The division is based on the attribute values of the data points. Unsupervised classification allows us to divide the dataset into several groups without knowing how the records "should" relate to each other. Sometimes the ways of grouping the records are very surprising and such study can bring many fresh ideas about the problems appeared in the area where data have been collected, and how to solve these problems.

The notion of a cluster is rather informal and we can give a formal definition of a cluster only in the framework of a formal mathematical model. Different formalizations are possible. In this paper we consider one of the most popular models, which we call the *point-based clustering model*.

An important characteristic of supervised and unsupervised classification methods is their *accuracy*. There are different approaches to *supervised classification accuracy* that can be used for comparing different classifiers (see for example, [10]). We mention here only *n*-fold cross-validation. This approach is very popular, however it cannot always give a good comparison of classifiers (see [3] for discussion). Indeed, the classification accuracy obtained for the same dataset in the case of  $n_1$ -fold cross-validation and  $n_2$ -fold crossvalidation are not necessarily the same if  $n_1 \neq n_2$ . The estimation of accuracy of clustering methods is much more difficult than supervised classification methods.

One of the approaches for evaluation of accuracy of clustering methods is based on the following idea: let us consider the dataset with given classes and apply a technique for clustering this dataset assuming that classes are unknown. Then the clusters obtained by this technique can be compared with classes and the accuracy of the technique can be estimated by the degree of coincidences of classes and clusters. This approach has been discussed in [7] where the notion of cluster "purity" was introduced (see also [6]). In this paper we show that comparison of classes and clusters is not always appropriate, hence this comparison cannot be always used for assessment of a clustering technique. Indeed, it is possible that the points have been grouped by this technique according to some other characteristic rather than the classes. A simple and interesting example of such a case can be found in [8]).

We now present another example. Assume that we have records that describe one of the two letters, either A or B typed in different fonts made of several typefaces (Arial, Times New Romans, Tahoma etc.), each written in two different styles: bold and italics. Different decomposition of this dataset onto two classes can be considered. For example, we can consider classes that contains fonts of A and fonts of B, respectively. We can also consider classes that contains both letters printed in bold and both letters printed in italics,

respectively. Different classes also can be considered. For example, one of them contains both letters typed in one collection of typeface and the other contains both letters typed in collection of different typefaces. Applying a clustering technique to this dataset we can divide it for two clusters based on a certain characteristic. For the sake of definiteness assume that one of these clusters contains letters printed in bold and the other contains letters printed in italics. On the other hand we can consider this dataset as a dataset with two classes, one of them consists of prints of the letter A and the other consists of prints of the letter B. In such a case the clusters and classes are pretty different, however we could not say the clustering technique is not good.

In such a situation we can suggest that some characteristics link points more strongly than their belonging to classes.

In this paper we show that a similar situation can appear in real-world datasets. We consider two real-world datasets with classes ("Pendigits" and "Letters") and compare classes and clusters for these datasets. We also give a very short description of results obtained in [13] for the Australian credit dataset that indicate another reason for non-coincidence between clusters and classes.

As was mentioned, we examine clusters in the framework of the pointbased clustering model. In this model we can use different optimization techniques for the search for clusters. We also apply the notions of cluster function and the notion of structure of clusters (see [6]) in order to check the quality of clusters obtained by this technique. The goal of this investigation is two-fold. First, we show that the relation between clusters and classes are very different in these datasets. Then we also compare different optimization techniques for real world datasets.

The paper is organized as follows. In section 2 we present a point-based clustering model and give a short description of non-smooth optimisation methods used in the implementation of this model. In section 3 we discuss the relations between clusters and classes for the Pendigits dataset where clusters are close to classes. In section 4 we discuss the relations between clusters and classes for the Letters dataset where classes and clusters are different. Section 5 summarizes the obtained results and draws some guidelines for future research.

## 2 Preliminaries

### 2.1 Clusters

The following definition can be found in [8]: Cluster Analysis is the organization of a collection of patterns (usually represented as a vector of measurements, or a point in a multidimensional space) into clusters based on similarity. Intuitively, patterns within a valid cluster are more similar to each other that they are to a pattern belonging to a different cluster.

Generally speaking, it is impossible to say whether given groups of points form clusters based on this definition (indeed, we need to explain what "more similar to each other" and "intuitively" mean). In order to find clusters we should formalize this definition. It can be done in many ways. We use one of the most popular models that we call a *point-based clustering model*. In the framework of this model clusters represent groups of points, centred in their centres, which are also points.

Suppose that a dataset contains N points  $a_i \in \mathbb{R}^n$ , i = 1, ..., N and suppose that we need to find k clusters in this dataset. In order to determine the clusters we have to find k points which are the centres of these clusters and then we assign each point to the cluster with the nearest centre. The search for the centres can be reduced to the minimization of the so-called *cluster* function:

$$f(x_1, \dots, x_k) = \sum_{i=1}^{N} \min_{j=1,\dots,k} \|a_i - x_j\|_p, \quad p \ge 1.$$
(1)

Thus we need to solve the following mathematical programming problem

minimize 
$$\sum_{i=1}^{N} \min_{j=1,\dots,k} \|a_i - x_j\|_p, \quad p \ge 1,$$
 (2)

subject to  $(x_1, \dots, x_k) \in \mathbb{R}^{nk}$  (3)

The dimension of this problem is nk (the number features multiplied by the number of clusters). If k > 1, the cluster function is nonsmooth and nonconvex. This function has many local minima, especially if the number of records in the dataset is large.

**Remark 1** There are some techniques which allow one to reduce the number of components in the sum (see [6]).

After finding centres, we construct clusters: each point is assigned to the cluster with the nearest centre.

The detailed explanation of the approach to clustering based on cluster function can be found in [6]. In the following we assume that p = 2.

**Remark 2** Models that are different from the point-based one also can be used. For example, clusters that are centred by some hyperplanes ("hyperplane-based clusters" or "skeletons") can be considered (see [14] and references therein).

Often we have different candidate vectors  $X = (x_1, \ldots, x_k)$  which contend to be considered as centres of clusters. In the framework of the point-based clustering model, the cluster function can be used as a tool for comparison of candidates. Namely, the vector  $X^* = (x_1^*, \ldots, x_k^*)$  is better suited for the role of centres of clusters then  $X = (x_1, \ldots, x_k)$  if  $f(x_1^*, \ldots, x_k^*) < f(x_1, \ldots, x_k)$ . See [6] for details.

## 2.2 Step by step clustering and simultaneous clustering

Finding an initial point for minimizing the cluster function is a problem whose difficulty grows with the number of clusters sought. For one cluster, the optimization problem is convex, however for  $k \ge 2$ , it is not anymore convex. A heuristic, step by step method was proposed in [4] for minimizing the cluster function.

This algorithm, which constructs the clusters one by one and refines the solutions stepwise, allows one to avoid the difficult task of finding a suitable initial point for the local method. This algorithm has been introduced and discussed in [4], where numerical results were also presented.

We can also try to find simultaneously all clusters if we have some initial guess. For example, comparing clusters and classes we can use the collection of centres of classes as an initial point for the optimization.

The k-means method is one of the most popular clustering methods (see, for example, [9], [11]). This clustering method is very fast, however, very often it produces clustering systems with higher dissimilarity than methods based on nonsmooth optimization technique that we use (see subsection 2.3). It is known that the results obtained by k-means highly depends on the initial points. It is efficient to use the k-means method to find a good starting point for nonsmooth optimisation (see [6] for further discussions).

### 2.3 Optimisation techniques

In our numerical experiments we use the discrete gradient method (briefly DG) for minimizing (2). DG has been developed by Adil Bagirov (see [1], [2] for details).

We also use a combination of DG with the cutting angle method (briefly CAM, see [5], [12] for details). We call this combination DG+CAM. The results obtained by DG and especially by DG + CAM do not depend so strongly on the initial point as those obtained by k-means.

DG and DG+CAM are the main optimisation techniques, provided by CIAO-GO (optimisation software developed at the University of Ballarat). This software has been developed for constrained and unconstrained optimisation. It allows researchers to choose different optimisation methods (DG, DG+CAM), different types of penalty functions (for constrained optimisation) and different types of initial points (initial points can be appointer by researchers or generated randomly in a feasible region). For more details about CIAO-GO software, refer to [15]

#### 2.4 Structure of clusters

As was mentioned in subsection 2.1, the cluster function can be used for evaluation of the quality of centers of clusters and through this for evaluation of quality of clusters. Another way to evaluate the quality of unsupervised classification is to check the distribution of the points within the clusters (structure of clusters). The goal is to check how "deep" the points are in the clusters.

Suppose that we work with a dataset which contains N observations. A clustering method has been applied to this dataset and k centres of clusters have been obtained. Consider a point a from the dataset which belongs to the *l*-th cluster (with the centre  $x^l$ ). For this point we determine a value c(a) which can be found as follows:

$$c(a) = \frac{||a - x^{l}||}{\min_{j=1,\dots,k, j \neq l} ||a - x^{j}||}.$$
(4)

From the definition of clusters it is easy to conclude that  $c(a) \in [0, 1]$ . Very often in our study we use a grid for c(a) such as  $0.1, 0.2, \ldots, 1$ . We do not

work with the exact value for each c(a), but use intervals. For example,

$$c(a) \in [(i-1) * 10^{-1}, i * 10^{-1}), \quad i = 1, \dots, 9 \text{ and } c(a) \in [0.9, 1]$$
 (5)

or

$$c(a) \in [0, i * 10^{-1})$$
  $i = 1, \dots, 9$  and  $c(a) \in [0, 1].$  (6)

We use the value c(a) for each point to describe how "deep" this point is inside the cluster. We should underline that different values of c(a) do not represent the radiuses for some spheres centred at the centres of the corresponding clusters. They rather represent some levels of confidence that the chosen point belongs to this cluster but not another one. It is possible that some points which are not "deep" enough inside the corresponding cluster move to another cluster (change their membership). It could happen, for example,

- if we change the norm in the definition of the cluster function;
- if we change the location of points (another accuracy to represent numbers in the computer);
- if we change the value for some internal parameters for the optimization methods etc.

If c(a) = 1 or close to 1 there are two centres such that the distances between the point and these two centres are (almost) the same. In this case some changes within the data may change the membership of the point (the point is "unstable" inside the cluster). If c(a) is close to 0 the level of confidence for the point to keep its membership is high (the point is "stable").

Suppose that we obtain two different clustering results. First, we can compare the value of cluster function in the centres of corresponding clusters. The systems of clusters with the lowest value of cluster function is better from the point of view of cluster function. Second, we check the structure of corresponding clusters by means of (4). If for the first of the clustering result the values c(a) are smaller than for the second clustering result for most of the points a, we assume that the first collection of the centres is preferable to the second one in the sense of the structure of the clusters (by means of (4)). The most important is to investigate the points a with the values  $c(a) \in [0.9, 1]$ . It is possible that the two approaches: evaluation of quality of clusters by means of clusters function and the structure of clusters, respectively, contradict each other (multi criteria problem).

## 2.5 Classes and clusters: purity

The notion of *purity* (see for example [7]) is used in the literature for evaluation of accuracy of clustering methods. Assume that we have a dataset A composed of the classes  $\{D_1, \ldots, D_l\}$  and we apply a clustering procedure for finding clusters  $\{C_1, \ldots, C_k\}$  in this dataset.

The purity of a set of clusters  $\{C_1, \ldots, C_k\}$  is calculated as follows:

$$p(\{C_1, \dots, C_k\}) = \frac{\sum_{i=1}^k \max_{j=1,\dots,l} |C_i \cap D_j|}{|A|},$$

where |B| is the cardinality of a finite set B.

We illustrate this notion in the simplest case where a dataset A contains only 2 classes  $\{D_1, D_2\}$ . Suppose that 4 clusters  $\{C_1, C_2, C_3, C_4\}$  have been found in this dataset.

If the majority of points from the j-th cluster (j = 1, ..., k) is in the i-th class (i = 1, 2), we assign the whole j-th cluster to the i-th class. When all the clusters are assigned to one of the classes, the percentage of the correctly classified points for the test set is considered as the classification accuracy.

$$p(\{C_1, \dots, C_4\}) = \frac{\sum_{i=1}^4 \max(|C_i \cap D_1|, |C_i \cap D_2|)}{|A|}$$

#### 2.6 Classes and clusters: preliminary example

There are many methods for supervised classification (see [10] and references within) which are based on quite different techniques. It is suggested, that if there is a group of points in a dataset which are misclassified by several methods then this group of points should be studied separately.

In this subsection we present some of the results obtained in [13] for comparing the division into classes and into clusters and for discovering possible reasons why sometimes these divisions are not the same.

The following procedure has been suggested in [13] for examination of quality of a dataset. First, different classifiers that are based on completely different ideas should be chosen. Then these methods should be applied to the dataset and all misclassified points for each method should be identified. The points that are misclassified by all methods are called *questionable*. Since the used methods are based on completely different approaches we can accept that misclassification of a questionable point does not depend on classifier and depends on the quality of the dataset.

It was indicated in [13] that the Australian credit dataset consists of 5 well-defined clusters. One of the clusters contains 56 points from the first class and 77 points from the second class. The majority of questionable points belongs to this cluster and they are "deep" inside the cluster.

The fact that questionable points are mainly placed in the same cluster is of great interest. One of the possible interpretation of this fact is as the following: questionable points appear as the result of a certain systematic error and they have no real links with both classes.

## 3 Pendigits

This dataset was introduced by E. Alpaydin and Fevzi Alimoglu (see [10]). It contains 10 classes, 10992 observations, 16 attributes. All input attributes are integers  $0 \dots 9$ .

The dataset has been created by collecting 250 samples from 44 writers. The samples written by 30 writers are used for the training set and the digits written by the other 14 are used for writer independent testing.

#### 3.1 Classes and centres

The experiment are first carried out on the Pendigits dataset. Although this dataset is quite large, it is usually well handled by most classification methods. The first experiments is to carry out a step by step clustering without any knowledge of the classes and compare these clusters with the classes. Table 1 shows the amount of points from each class in each cluster.

As a result it is noticeable that a great majority of points in each cluster belongs to the same class. This means that the clusters seem to coincide with classes. The purity is 78.58%.

The second experiment carried out is to find for each class one centre. The points obtained are then considered as cluster centres and the same analysis as previously is applied. Table 2 presents the results.

A similar - but more previsible - observation is made about this table: the correspondence between clusters and classes is very strong. The purity is

	1	2	3	4	5	6	7	8	9	10	size	% class
1	470	0	0	0	0	107	2	1	1	0	581	80.89
2	1	1122	334	5	1	1	0	0	6	1	1471	76.27
3	0	15	634	26	0	13	0	123	146	31	988	64.17
4	0	0	2	1051	5	28	0	79	0	1	1166	90.13
5	11	0	1	29	1049	3	1	0	1	0	1095	95.79
6	247	0	0	0	0	989	0	12	0	0	1248	79.24
7	183	0	0	0	1	0	625	0	3	0	812	76.97
8	8	0	82	32	0	2	191	715	0	2	1032	69.28
9	69	2	1	0	0	0	0	0	962	0	1034	93.03
10	66	5	89	1	0	0	236	125	23	1020	1565	65.17
size	1055	1144	1143	1144	1056	1143	1055	1055	1142	1055		
	44.54	98.07	55.46	91.87	99.33	86.52	59.24	67.77	84.23	96.68		
												-

Table 1

Pendigits: repartition of the classes in the clusters obtained by step by step clustering

	1	2	3	4	5	6	7	8	9	10	size	% class
1	561	0	0	0	0	142	14	1	0	0	718	78.13
2	0	1101	325	9	2	1	0	0	7	1	1446	76.14
3	0	37	651	18	0	12	1	126	148	31	1024	63.57
4	0	0	3	1060	9	79	0	34	0	1	1186	89.37
5	42	0	0	11	1043	25	3	0	0	0	1124	92.79
6	157	0	0	0	0	876	0	8	0	0	1041	84.14
7	67	0	0	0	1	0	613	0	0	0	681	90.01
8	6	0	75	45	0	8	180	770	0	2	1086	70.9
9	154	1	0	0	0	0	0	0	968	0	1123	86.19
10	68	5	89	1	1	0	244	116	19	1020	1563	65.25
size	1055	1144	1143	1144	1056	1143	1055	1055	1142	1055		
	53.17	96.24	56.95	92.65	98.76	76.64	58.1	72.98	84.76	96.68		

Table 2

Pendigitits: repartition of the classes in the clusters defined by class centres

	0	.1	.2	.3	.4	.5	.6	.7	.8	.9	1	mean
cluster 1	L											0.58
cluster 2	2											0.52
cluster 3	3											0.52
cluster 4	1											0.43
cluster 5	5											0.56
cluster 6	3											0.51
cluster 7	7											0.49
cluster 8	3											0.65
cluster 9	)											0.59
cluster 10	)											0.56

Fig. 1. Structures of the clusters of Pendigits

78.81%.

The class centres are used as an initial point to the local optimisation method to find a local minimum to the cluster function. In the case of Pendigits the solution obtained is the same as the one reached by the step by step method.

Because the same result was reached by two different methods, it can be expected that it is a very good local minimum. For comparison, the objective function value for a solution reached by the k-means method with a randomly chosen initial point was more than 2% larger.

The conclusion of this is that even by applying a clustering method without any knowledge of classes, the solution obtained is the one which corresponds at best to the classes. The classification accuracy is relatively high. In the Pendigits dataset the classes and cluster seem to be perfectly equivalent.

Indeed when the number of clusters is increased to 20 (twice the number of classes) for the step by step clustering, the purity becomes 87%. This result, obtained by a method which does not use any knowledge of classes, is comparable to most specifically designed methods presented in [10].

## 3.2 Classes and cluster structures

Figure 1 represents the structure of the clusters, The darker the square the larger amount of points are present in the layer. All the clusters present a similar structure: the points are deep, showing a good clusterization.

The Pendigits dataset has been shown to present a strong correlation between clusters and classes. It may be interesting to consider a deeper relation, by finding the layers of the clusters.



Fig. 2. Structures of the clusters in Pendigits

Figure 2 shows the depth of the classes in each cluster. Each disk represents one cluster, and each slice of this disk represents one class. The darker one circle is, the larger the proportion of points from this class is deeper inside the cluster.

Table 3 shows the repartition of the main class and the other points in each cluster. All the clusters present the same characteristic: although the points of each clusters are quite deep for all the clusters the main class is deeper than the others. This means that not only the clusters centres represent the classes well, but also the "misclassified points" do not belong so strongly to the cluster.

Table 4 shows the average depth of each class inside each cluster. The results presented in this table confirm that in each cluster the majoritary class is the deepest.

### 3.3 Summary

- The relationships between classes and clusters are very explicit
- The points are "deep" inside the clusters, and therefore the clustering is considered of high quality
- The correctly classified points are deeper.

## 4 Letters

The dataset was introduced by David Stale (see [10]). It is based on various fonts representation of the letters of the Latin alphabet. The dataset consists of 20000 observations, 26 classes, 16 numerical attributes. This is

Cluster		[0, 0.1]	[0.1, 0.2]	[0.2, 0.3]	[0.3, 0.4]	[0.4, 0.5]	[0.5, 0.6]	[0.6, 0.7]	[0.7, 0.8]	[0.8, 0.9]	[0.9,1]
1	Main class	0	0	6.8	22.12	18.72	16.17	11.06	11.27	6.17	0
	Other classes	0	0	0	0	0	3.6	37.83	18.01	24.32	0
2	Main class	0	0.08	5.97	22.01	26.73	18.36	13.72	8.55	3.83	0
	Other classes	0	0	1.43	14.61	27.79	16.61	15.47	9.74	7.73	0
3	Main class	0.31	11.35	18.61	23.34	17.5	12.77	6.46	5.36	2.52	0.31
	Other classes	0	0	2.25	6.49	16.94	12.42	6.49	15.53	19.2	0
4	Main class	0	1.23	10.56	17.6	21.4	13.32	14.93	8.65	7.61	0
	Other classes	0	0	0	0	0	0	0.86	12.17	38.26	0
5	Main class	0	5.71	22.11	25.64	20.59	13.34	7.14	3.62	1.42	0
	Other classes	0	0	0	0	0	0	0	10.86	19.56	0
6	Main class	0	4.24	23.55	19.61	18.7	12.84	8.08	5.66	3.74	0
	Other classes	0	0	0	0	0.77	3.47	10.42	22.39	33.59	0
7	Main class	0	7.68	26.08	26.24	19.36	13.44	4.32	1.76	0.96	0
	Other classes	0	0	0.53	0	0	1.6	6.95	16.57	28.34	0
8	Main class	0	0	0.97	11.74	20.27	17.34	13.28	10.48	14.82	0
	Other classes	0	0	0.31	3.47	10.72	13.56	14.51	26.81	16.71	0
9	Main class	0	0.1	4.67	17.15	23.38	22.03	14.55	9.04	5.5	0
	Other classes	0	0	0	0	0	0	1.38	15.27	36.11	0
10	Main class	0	0.49	13.23	26.17	23.62	18.82	8.43	4.7	3.13	0
	Other classes	0	0	0	0.18	1.28	6.42	14.67	22.2	26.42	0

Table 3

Pendigits: repartition of the classes by cluster layers

samples of 26 capital letters, printed in different fonts. 20 fonts have been considered and the location of the corresponding samples has been distributed randomly within the dataset.

The experiments described and applied on the Pendigits dataset in the previous section are here applied on the Letters dataset. The results are presented and discussed.

	1	2	3	4	5	6	7	8	9	10
1	0.55	-	-	-	-	0.76	0.77	0.87	0.91	-
2	0.99	0.51	0.56	0.91	0.98	0.88	-	-	0.96	0.56
3	-	0.81	0.41	0.84	-	0.88	-	0.83	0.5	0.78
4	-	-	0.85	0.53	0.93	0.93	-	0.87	-	0.85
5	0.92	-	0.91	0.91	0.41	0.99	0.97	-	0.99	-
6	0.82	-	-	-	-	0.45	-	0.89	-	-
7	0.86	-	-	-	0.26	-	0.38	-	0.88	-
8	0.93	-	0.81	0.86	-	0.97	0.62	0.63	-	0.76
9	0.87	0.96	0.99	_	-	_	_	_	0.54	-
10	0.86	0.9	0.77	0.98	-	-	0.76	0.86	0.94	0.47

Table  $\overline{4}$ 

Average depth for each class in each cluster

### 4.1 Classes and centres

The Letters dataset contains 26 classes. Therefore, 26 clusters are found using the step by step algorithm. The size of the intersection of each cluster with each class is then evaluated.

Table 5 presents the results.

Unlike for the Pendigits dataset, where each cluster is strongly associated with a class (and where each class is represented by one cluster), the repartition of the classes in the clusters of the Letters dataset is much more scattered. The purity is very low: 28.25%

Table 6 presents the repartition of the classes when the cluster centres considered are one centre per class. It is noticeable that here again the repartition of the points is very disseminate. The purity is 23.94%, the value of the cluster function is much higher than in the case of the step by step clustering. We suggest that the classes are not "point-based".

Finally, these class centres are used as an initial point for the local algorithm to minimize the cluster function. While in the case of the Pendigits dataset the result of this experiment was similar to the solution found by step by step clustering, this time the centres are quite far apart and another solution is reached. Once again, the repartition, shown in Table 7 is very diffuse.

A possible interpretation of these results could be that the clustering

	34	18	$\infty$	91	6	30	35	20	23	41	92	17	23	6	100	11	51	38	40	30	97	12	22	16	19	100			
	708	676	949	301	1298	886	485	555	191	703	616	880	1137	1315	324	1047	680	588	805	984	174	978	784	1469	1299	168			
26	0	0	6	3	57	9	0	0	0	0	2	0	15	26	0	12	190	190	-	2	0	12	0	2	52	168	747	25	50
25	0	95	47	0	83	0	0	33	0	0	0	14	33	111	0	98	4	0	0	0	0	0	0	23	255	0	766	33	erin
24	34	41	40	0	ъ	0	0	36	0	0	33	66	150	66	0	Ч	0	0	Ч	$\infty$	0	c,	0	$^{249}$	37	0	773	32	clust
23	104	97	21	0	0	68	0	ъ	0	0	0	21	70	34	0	0	0	0	$\infty$	17	0	87	40	123	41	0	736	16	ep (
22	0	0	10	0	15	157	0	5	0	198	0	0	14	50	0	2	0	0	0	с	0	127	174	ъ	39	0	262	$^{24}$	oy st
21	0	0	59	0	0	105	0	35	37	2	0	66	45	20	0	0	0	0	26	41	169	126	28	54	0	0	813	20	ep þ
20	0	0	0	0	0	0	0	26	25	0	0	0	0	18	0	0	0	0	327	304	0	0	2	0	0	0	752	43	y st
19	0	0	15	ŝ	74	-	0	0	0	0	Ч	0	0	16	0	30	350	229	0	-	0	2	0	4	29	0	755	46	ed b
18	0	0	44	0	-	125	0	16	0	212	0	0	4	26	0	0	0	0	Ч	-	0	119	146	2	36	0	786	26	tain
17	1	81	39	0	111	11	0	0	0	0	0	45	×	55	0	122	0	Ч	0	0	0	69	4	128	112	0	787	16	qo s
16	0	0	25	0	73	0	0	5 L	0	0	5	0	13	34	324	ъ	00	158	Ч	-	0	2	0	6	46	0	192	42	sters
15	0	0	32	0	43	0	0	18	0	0	4	134	138	120	0	-	0	0	0	2	0	0	0	187	104	0	783	23	clu
14	0	0	53	0	16	0	0	29	0	0	0	153	266	11	0	0	0	0	9	4	0	0	0	210	5	0	753	35	$_{\mathrm{the}}$
13	187	107	15	0	79	4	0	9	0	0	0	0	23	66	0	116	Ч	0	0	1	0	44	3 S	16	67	0	768	24	s in
12	0	35	71	0	74	×	0	30	14	0	Ч	32	35	76	0	93	0	0	-	20	0	92	6	119	29	0	739	16	lasse
11	0	0	36	0	28	0	0	Ч	0	0	567	0	0	44	0	-	59	0	13	21	0	1	0	-	17	0	789	71	le c
10	0	0	34	2	0	91	0	$^{24}$	0	291	0	0	Ч	45	0	0	0	0	26	27	0	84	137	0	5	0	764	38	of t]
6	0	0	67	0	111	0	0	4	0	0	ŝ	93	169	65	0	114	4	0	18	6	0	0	0	123	25	0	805	20	ion
8	0	0	2	0	0	0	159	116	44	0	×	0	0	13	0	0	0	0	239	203	0	2	2	0	4	0	792	30	artit
4	0	4	62	0	103	0	0	31	0	0	0	88	85	80	0	89	-	0	0	11	0	0	0	148	53	0	758	19	rep
9	0	0	40	0	11	2	171	27	40	0	Ч	36	27	27	0	1	0	0	127	251	0	7	ъ	6	Н	0	783	32	ers:
ъ	0	0	27	17	61	274	0	°°	0	0	0	0	0	52	0	33	4	5	×	10	0	79	163	0	39	0	775	35	Lett
4	0	0	69	0	63	3 S	155	$^{24}$	31	0	ŝ	43	40	69	0	72	0	0	2	34	ъ	40	3 S	50	28	0	734	21	5.
3	140	06	32	0	125	4	0	9	0	0	15	0	0	52	0	98	0	0	0	0	0	41	0	0	142	0	748	18	able
~	0	0	84	276	50	21	0	28	0	0	0	56	28	66	0	55	4	ъ	0	10	0	26	68	4	16	0	803	34	L
					2	~	_	_	C	0	ŝ	0	0	23	0	-04	0	0	0	0	0	15	0	0	93	0	34	32	
1	242	123	16	0	11	0	0	0	-							_											1~		

123456111		39	0	100	23	34	26	20	27	100	50	09	45	24	50	92	22	15	11	21	29	4	24	53	22	36	25			
123456789101112131617181920212223242524252425242524252425242524242424111121000 <t< td=""><td></td><td>283</td><td>0</td><td>242</td><td>1179</td><td>925</td><td>30</td><td>375</td><td>441</td><td>ŝ</td><td>2</td><td>46</td><td>768</td><td>178</td><td>908</td><td>56</td><td>22</td><td>281</td><td>4265</td><td>1205</td><td>183</td><td>2612</td><td>2909</td><td>358</td><td>1577</td><td>767</td><td>385</td><td></td><td></td><td></td></t<>		283	0	242	1179	925	30	375	441	ŝ	2	46	768	178	908	56	22	281	4265	1205	183	2612	2909	358	1577	767	385			
1234567891011121316171810 </td <td>26</td> <td>16</td> <td>0</td> <td>0</td> <td>22</td> <td>16</td> <td>Ч</td> <td>0</td> <td>19</td> <td>0</td> <td>0</td> <td>0</td> <td>0</td> <td>0</td> <td>64</td> <td>0</td> <td>0</td> <td>4</td> <td>54</td> <td>77</td> <td>31</td> <td>89</td> <td>6</td> <td>44</td> <td>225</td> <td>76</td> <td>0</td> <td>747</td> <td>30</td> <td></td>	26	16	0	0	22	16	Ч	0	19	0	0	0	0	0	64	0	0	4	54	77	31	89	6	44	225	76	0	747	30	
1         2         3         4         5         6         7         8         9         10         1	25	4	0	0	179	0	0	54	13	0	0	0	0	5	Ч	0	0	5	247	25	Ч	198	37	0	0	0	0	, 992	32	
12345678111314151617181920212223111111000<	$^{24}$	e S	0	0	14	0	0	0	4	0	0	0	0	0	4	0	0	30	372	x	0	110	41	0	0	83	98	, 222	48	res
1234567811131415161718192021231113150000001000 <t< td=""><td>23</td><td>9</td><td>0</td><td>0</td><td><math>\infty</math></td><td>0</td><td>0</td><td>0</td><td>0</td><td>0</td><td>0</td><td>0</td><td>2</td><td>0</td><td>0</td><td>0</td><td>0</td><td>П</td><td>215</td><td>0</td><td>0</td><td>63</td><td>42</td><td>46</td><td>66</td><td>282</td><td>ъ</td><td>736</td><td>38</td><td>cent</td></t<>	23	9	0	0	$\infty$	0	0	0	0	0	0	0	2	0	0	0	0	П	215	0	0	63	42	46	66	282	ъ	736	38	cent
11234567891011121316171819202011131500 <t< td=""><td>22</td><td></td><td>0</td><td>0</td><td><math>\infty</math></td><td>194</td><td>0</td><td>0</td><td>0</td><td>0</td><td>0</td><td></td><td>86</td><td>0</td><td>0</td><td>0</td><td>0</td><td>0</td><td>82</td><td>0</td><td>0</td><td>54</td><td>4</td><td>0</td><td>358</td><td>0</td><td>ъ</td><td>962</td><td>44</td><td>ass</td></t<>	22		0	0	$\infty$	194	0	0	0	0	0		86	0	0	0	0	0	82	0	0	54	4	0	358	0	ъ	962	44	ass
12345678910111213161718192011131500000000000000002000000000000000000300000000000000000420000000000000000061100000000000000000611000 <t< td=""><td>21</td><td>0</td><td>0</td><td>0</td><td>0</td><td>0</td><td>0</td><td>0</td><td>11</td><td>0</td><td>0</td><td>0</td><td>22</td><td>0</td><td>2</td><td>0</td><td>0</td><td>23</td><td>399</td><td>0</td><td>0</td><td>13</td><td>126</td><td>190</td><td>5</td><td>0</td><td>25</td><td>313 '</td><td>49</td><td>y cl</td></t<>	21	0	0	0	0	0	0	0	11	0	0	0	22	0	2	0	0	23	399	0	0	13	126	190	5	0	25	313 '	49	y cl
12345678910111316171819111315009000100000020000000000000000210000000000000000310111300000000000000420227000000000000004200000000000000000617100000000000000000617100 </td <td>20</td> <td>0</td> <td>32</td> <td>3</td> <td>0</td> <td>0</td> <td>717</td> <td>0</td> <td>0</td> <td>0</td> <td>0</td> <td>752 8</td> <td>95</td> <td>ed b</td>	20	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	32	3	0	0	717	0	0	0	0	752 8	95	ed b
11234567891011121316161718111315000000000000002000000000000000003112000000000000000420000000000000000042000000000000000006110110000000000000070000000000000000070000000000000000000700000000000000000010000000000 <td>19</td> <td>12</td> <td>0</td> <td>0</td> <td>11</td> <td>14</td> <td>0</td> <td>ъ</td> <td>23</td> <td>0</td> <td>0</td> <td>0</td> <td>9</td> <td>35</td> <td>37</td> <td>0</td> <td>0</td> <td>0</td> <td>ŝ</td> <td>138</td> <td>15</td> <td>54</td> <td>` ຕ</td> <td>0</td> <td>347</td> <td>4</td> <td>48</td> <td>755 '</td> <td>45</td> <td>lefin</td>	19	12	0	0	11	14	0	ъ	23	0	0	0	9	35	37	0	0	0	ŝ	138	15	54	` ຕ	0	347	4	48	755 '	45	lefin
11345678910111314151617111130000000000000020000000000000000030100 <td>18</td> <td>0</td> <td>0</td> <td>0</td> <td></td> <td>127</td> <td>0</td> <td>0</td> <td>ъ</td> <td>0</td> <td>0</td> <td>2</td> <td>177</td> <td>28</td> <td>0</td> <td>0</td> <td>0</td> <td>0</td> <td>83</td> <td>0</td> <td>0</td> <td>118</td> <td>9</td> <td>-</td> <td>201</td> <td>0</td> <td>37</td> <td>786</td> <td>25</td> <td>ers c</td>	18	0	0	0		127	0	0	ъ	0	0	2	177	28	0	0	0	0	83	0	0	118	9	-	201	0	37	786	25	ers c
1         2         3         4         5         6         7         8         9         10         11         12         34         15         6         7         8         9         10         11         12         14         15         0         9         0         0         0         16         15         16         15           2         0	17	0	0	0	165	16 ]	1	34	30	0	0		51 ]	2	131	0	0	0	27	43	54	507	6	5 L	9	0	5 L	1 287	26	uste
12345678910111314151113150090000000020000000000000302420000000000004220227378000000000651111017000000000065111017000000000061011<	16	15	0	0	17	0	4	2	57	0	0	0	0	0	92 ]	0	0	0	53	262	16	74 2	18	0	147		0	2 192	34	te cl
1234567891011121314111315009000036020000000000360302420000000000004220227378032010360000511110130320100000000651110172807800000060000000000000060000000000000060000000000000060000000000000600000000000001000000000000011000<	15	16	0	0	19	0	0	4	11	0	0	0	0	0	ŝ	0	0	43	<b>1</b> 73	41	0	115	11	0	0	0	47	783 7	60	in tł
123456789101112131113150000000000200000000000002000000000000003024200000000000042202710000000000060000000000000006000000000000000733500000000000006000000000000000600000000000000700000000000000100000000000 <td< td=""><td>14</td><td>0</td><td>0</td><td>0</td><td>0</td><td>0</td><td>0</td><td>0</td><td>6</td><td>0</td><td>0</td><td>0</td><td>0</td><td>0</td><td>2</td><td>0</td><td>0</td><td>44</td><td>157 4</td><td>96</td><td>0</td><td>20</td><td>48</td><td>5</td><td>0</td><td>0</td><td>19</td><td>753</td><td>60</td><td>ses</td></td<>	14	0	0	0	0	0	0	0	6	0	0	0	0	0	2	0	0	44	157 4	96	0	20	48	5	0	0	19	753	60	ses
123456789101112111315000000000200000000000030242000000000004220227378084071000671121130000000006171121302420000000600000000000073355111017280000060000000000007335511101728000010000000000000111112121312131414141413150000000014141500000000 </td <td>13</td> <td>36</td> <td>0</td> <td>0</td> <td>77</td> <td>18</td> <td>0</td> <td>36</td> <td>1</td> <td>0</td> <td>0</td> <td>0</td> <td>5</td> <td>37</td> <td>0</td> <td>0</td> <td>0</td> <td>0</td> <td>134</td> <td>Ч</td> <td>10</td> <td>205</td> <td>0</td> <td>×</td> <td>ъ</td> <td>195</td> <td>0</td> <td>1 892</td> <td>26</td> <td><math>_{\rm clas}</math></td>	13	36	0	0	77	18	0	36	1	0	0	0	5	37	0	0	0	0	134	Ч	10	205	0	×	ъ	195	0	1 892	26	$_{\rm clas}$
123456789101111131500000100021150000000000302420000000000042202273780840700065000000000006500000000000600000000000060000000000006000000000000600000000000060000000000000600000000000001000000000000011000000 <td< td=""><td>12</td><td>0</td><td>0</td><td>0</td><td>59</td><td>0</td><td>0</td><td>5 L</td><td>2</td><td>0</td><td>-</td><td>0</td><td>44</td><td>c,</td><td>0</td><td>52</td><td>17</td><td>П</td><td>28</td><td>°</td><td>50</td><td>42 2</td><td>81</td><td>19</td><td>0</td><td>0</td><td>32</td><td>39 7</td><td>24</td><td><math>_{\mathrm{the}}</math></td></td<>	12	0	0	0	59	0	0	5 L	2	0	-	0	44	c,	0	52	17	П	28	°	50	42 2	81	19	0	0	32	39 7	24	$_{\mathrm{the}}$
12345678910111315000001021000000000302420000000004200227378024000005171121303201000000650780320100000733554270000000083305111017280034611070000000011011101728003461101011101728140121415000000113150116101610101141501101110111150110101010101160000010101<	11	0	0	0	0	0	0	0	16	0	0	0	0	0	155	0	0	0	28 ]	136	0	39 ]	[15]	0	0	0	0	1 682	57	1 of
12345678911131500000002000000000030242000000004220227378084000651711213032010000651711213032010000651711100780006517111000171012733554171717280001100000000000110111000000001101110111011101113121314101110111313131500000000141501110111313131516131611101413131613<	10	0	0	0	0	77	0	0		3	-	28	346	44	1	0	5 L	П	28	0	0	40	110	4	15	0	60	764	45	itio
1         2         3         4         5         6         7         8           1         113         1         50         0         9         0         0         0           2         0         0         0         0         0         0         0         0         0           3         0         242         0         7         8         0         9         0 <td>6</td> <td>1</td> <td>0</td> <td>0</td> <td>5</td> <td>0</td> <td>0</td> <td>38</td> <td>123</td> <td>0</td> <td>0</td> <td>0</td> <td>0</td> <td>0</td> <td>42</td> <td>0</td> <td>0</td> <td>5</td> <td>312</td> <td>[]]</td> <td>0</td> <td>139</td> <td>32</td> <td>0</td> <td>0</td> <td>0</td> <td>0</td> <td>305</td> <td>38</td> <td>part</td>	6	1	0	0	5	0	0	38	123	0	0	0	0	0	42	0	0	5	312	[]]	0	139	32	0	0	0	0	305	38	part
1         2         3         4         5         6         7           1         113         1         50         0         0         0         0           2         0         0         0         0         0         0         0         0           3         0         242         0         7         8         0         84           4         220         2         33         7         8         0         84           6         5         11         13         0         320         1         0           6         5         1         10         10         17         28           7         33         5         54         27         10         17         28           7         33         5         11         10         17         28         29         10         10           10         0         7         11         10         17         28         10         10         10           11         0         1         10         11         10         17         28         40           11 <t< td=""><td>8</td><td>0</td><td>0</td><td>0</td><td>0</td><td>0</td><td>0</td><td>0</td><td>0</td><td>0</td><td>0</td><td>0</td><td>0</td><td>0</td><td>0</td><td>0</td><td>0</td><td>15</td><td>188</td><td>0</td><td>0</td><td>-</td><td>586</td><td>2</td><td>0</td><td>0</td><td>0</td><td>792 8</td><td>73</td><td>s: re</td></t<>	8	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	15	188	0	0	-	586	2	0	0	0	792 8	73	s: re
1         2         3         4         5         6           1         113         1         50         0         9         0           2         0         0         0         0         0         0         0           3         0         242         0         0         320         1           4         220         2         273         7         8         0           5         17         112         13         0         320         1           6         7         11         10         17         0         0           7         33         5         54         27         0         0           7         33         5         11         10         17         0         0           8         3         39         5         11         10         17         0         0           10         0         7         11         10         10         17         0         11           11         0         31         10         10         17         10         11           11         11         1	7	0	0	0	84	0	x	78	28	0	0	0	0	0	18	0	0	40	219	46	0	96	135	0	0	ъ		, 857	28	tter
1         2         3         4         5           1         113         1         50         0         9           2         0         0         0         0         0         0           3         0         242         0         0         0         0         0           4         220         2         273         7         8         320           6         5         112         13         0         320           7         33         5         54         27         0           7         33         5         54         27         0           8         33         5         54         27         0           9         0         0         0         0         0         0           10         0         5         11         10         0         0           11         0         5         31         6         0         0           11         0         5         31         6         0         0           11         0         5         31         6         0         0 </td <td>9</td> <td>0</td> <td>0</td> <td>0</td> <td>0</td> <td>Ч</td> <td>0</td> <td>0</td> <td>17</td> <td>0</td> <td>0</td> <td>9</td> <td>0</td> <td>0</td> <td>1</td> <td>4</td> <td>0</td> <td>28</td> <td>216</td> <td>25</td> <td>Ч</td> <td>4</td> <td>171</td> <td>9</td> <td>0</td> <td>0</td> <td>0</td> <td>, 887</td> <td>60</td> <td>. Le</td>	9	0	0	0	0	Ч	0	0	17	0	0	9	0	0	1	4	0	28	216	25	Ч	4	171	9	0	0	0	, 887	60	. Le
1         2         3         4           1         113         1         50         0           2         0         0         0         0         0           3         0         242         0         0         1           4         220         2         73         7           5         17         112         13         0         1           6         5         0         7         1         1           7         33         5         54         27         1           7         33         5         54         27         1           8         3         39         5         11         6         11           9         0         0         0         0         0         11         6         11           10         0         31         1         6         11         1         6         11           11         0         31         1         6         31         6         11         1         6         11         1         1         1         1         1         1         1	5	6	0	0	×	320	0	0	10	0	0	0	19	0	0	0	0	0	20 2	56	0	116	41	2	171	0	0	775	41	ole 6
$\begin{array}{ c c c c c c c c c c c c c c c c c c c$	4	0	0	0	2	0	Ч	27	11	0	0	0	9	4	9	0	0	26	299	39	4	157	120	$^{24}$	0	0	3	734	40	$\operatorname{Tab}$
$\begin{array}{ c c c c c c c c c c c c c c c c c c c$	3	50	0	0	273	13	2	54	ы	0	0	ŝ	Ч	ъ	31	0	0	0	28	10	0	194	0	2	10	12	0	748 '	36	
$\begin{array}{ c c c c c c c c c c c c c c c c c c c$	2		0	242	5	112	0	ъ	39	0	0	ъ	e	0	0	0	0	21	76	72	0	165	44	0	15	1	0	803 '	30	
1     1       1     1       1     1       2     5       1     1       1 <td>-</td> <td>113</td> <td>0</td> <td>0</td> <td>220</td> <td>17</td> <td>5 L</td> <td>33</td> <td>3</td> <td>0</td> <td>0</td> <td>0</td> <td>0</td> <td>15</td> <td>15</td> <td>0</td> <td>0</td> <td>0</td> <td>42</td> <td>13</td> <td>Ч</td> <td>140</td> <td>0</td> <td>0</td> <td>6</td> <td>108</td> <td>0</td> <td>734 8</td> <td>29</td> <td></td>	-	113	0	0	220	17	5 L	33	3	0	0	0	0	15	15	0	0	0	42	13	Ч	140	0	0	6	108	0	734 8	29	
		-	2	3	4	5	9	4	×	6	10	11	12	13	14	15	16	17	18	19	20	21	22	23	$^{24}$	25	26	size		

algorithm does not reach meaningful clusters. The number of clusters is fairly large, and the cluster function possess a very large number of local minima. Its minimization is a very difficult task, even for a powerful method like the step by step method.

However neither the improved class centres nor the class centres themselves seem to represent the classes well. Moreover the value of the function at the point reached by the step by step method is sensibly lower than the one at the other points.

The distribution of the points of the dataset in clusters independent from their class may be due to the fact that some more information is contained in the dataset (for example the font of the letters), and this information may bring some noise for an eventual classification.

The latter interpretation is confirmed by the fact that when the number of clusters is decreased to 4, thus making the problem easier to solve, the repartition of the classes remains scattered.

## 4.2 Classes and cluster structures

Figure 3 shows the structure of the clusters. These cluster present generally a similar structure: for most of the clusters the majority of points is deep inside the cluster. This shows that the clustering created clusters of good quality.

Experiments show that the distribution of the classes inside the clusters is very diffuse for the letters dataset. This means that the classes do not represent the most obvious clusterization of this dataset. Let us consider the structures of the clusters and examine the distributions of the classes among the layers.

Figure 4 and table 8 show the distribution of the classes in each cluster. Clearly there are several types of clusters. Some of them contain only points from a few classes. The majority of the clusters, however, contains points belonging to many classes. In general, it can be noticed that the repartition in layers for all the classes present in the cluster is quite similar. This result emphasizes the conclusion that a cluster cannot be associated with a particular class in the Letters dataset.

In this dataset 20 different fonts have been considered and the location of the corresponding samples has been distributed randomly within the dataset. The hypothesis is that the meaning of the clusters is font recognition rather than letter recognition. However, it is impossible to check this hypothesis,

	15	28	0	0	0	x	13	36	10	26	12	38	76	66	13	22	10	100	12	44	32	39	15	16	100	60			
	760	848	0	0	0	1032	817	953	977	688	804	666	475	301	975	1249	1607	326	1272	835	956	870	1329	1570	2	355			
26	0	0	0	0	0	16	13	2	5	0	9	0	15	0	0	14	67	0	63	285	Ч	2	41	4	2	214	747	38	res
25	109	94	0	0	0	45	78	0	74	0	84	0	0	0	10	23	134	0	62	0	0	0	0	53	0	0	766	17	cent
24	63	42	0	0	0	43	S	0	67	0	2	0	S	0	56	151	72	0	33	0	$\infty$	2	0	254	0	0	773	32	ass
23	108	85	0	0	0	31	Ч	71	18	0	Ч	0	0	0	22	77	58	0	4	0	17	11	103	129	0	0	736	17	le cl
22	5	0	0	0	0	23	25	131	12	0	Ч	300	0	0	0	16	53	0	19	0	5	0	205	4	0	0	206	37	ıg tł
21	0	0	0	0	0	63	6	58	49	169	4	9	0	0	107	38	ъ	0	0	0	40	26	171	68	0	0	813	21	ovit
20	0	0	0	0	0	0	0	0	62	0	0	0	0	0	15	0	0	0	0	0	313	343	2	0	0	0	752	45	mpr
19	5	0	0	0	0	26	12	7	Ч	0	6	0	7	0	45	0	38	0	88	374	Ч	0	9	3	0	141	755	49	by i
18	0	0	0	0	0	46	6	89	51	0	0	307	0	0	0	9	79	0	°	0	1	Ч	192	5	0	0	786	39	ned
17	92	4	0	0	0	44	51	2	13	0	80	0	0	0	40	13	101	0	154	0	0	0	61	129	0	0	787	19	defi
16	33	0	0	0	0	36	x	0	23	0		0	13		0	15	55	326	101	168	0	Ч	0	10	0	0	761	42	ers
15	0	3	0	0	0	45	20	0	49	0	0	0	4	0	125	130	176	0	36	0	2	Ч	0	192	0	0	783	24	elust
14	0	0	0	0	0	59	11	0	32	0	0	0	2	0	133	283	2	0	11	0	9	9	0	203	0	0	753	37	the c
13	107	192	0	0	0	15	57	1	26	0	100	0	0	0	0	19	108	0	103	0	1	0	22	17	0	0	768	25	in t
12	48	0	0	0	0	68	89	2	60	0	72	0	5	0	37	45	60	0	64	0	18	2	61	111	0	0	739	15	SSeS
11	0	0	0	0	0	25	21	0	11	0	Ч	0	364	299	0	0	38	0	0	0	20	6	0	-	0	0	789	46	e cla
10	0	0	0	0	0	41	0	26	42	0	0	385	0	0	0	0	25	0	0	0	22	32	190	-	0	0	764	50	f the
6	0	0	0	0	0	72	84	0	23	0	100	0	2	0	62	175	55	0	20	0	6	15	0	116	0	0	805	21	o uo
x	0	0	0	0	0	Ч	Ч	0	101	181	0	0	12	0	20	0	6	0	0	0	196	270	Ч	0	0	0	792	34	titic
4	15	0	0	0	0	61	71	0	65	0	73	0	0	П	69	66	73	0	75	0	12	0	0	144	0	0	758	18	epaı
9	0	0	0	0	0	36	10	0	41	178	3	Ч	13	0	72	25	11	0	0	0	238	137	10	x	0	0	783	30	rs: r
ഹ	0	0	0	0	0	36	35	344	21	0	18	0	0	0	22	0	66	0	49	ъ	10	9	162	Ч	0	0	775	44	ette
4	0	0	0	0	0	54	110	1	39	160	65	0	3	0	70	30	55	0	32	0	31	ъ	23	56	0	0	734	21	7. L
с.	67	183	0	0	0	37	24	1	26	0	76	0	25	0	0	0	111	0	135	0	0	0	27	9	0	0	748	24	ble
5	0	0	0	0	0	06	41	223	51	0	29	0	0	0	53	90	57	0	54	S	×	Ч	49	54	0	0	803	27	$\operatorname{Ta}$
-	114	$^{242}$	0	0	0	19	32	0	Ч	0	79	0	3	0	0	0	94	0	146	0	0	0	S	Ч	0	0	734	32	
	1	2	°	4	ъ	9	7	x	6	10	11	12	13	14	15	16	17	18	19	20	21	22	23	$^{24}$	25	26	size		

(	)	.1	.2	.3	.4	.5	.6	.7	.8	.9	1	mean
cluster $1$												0.80
cluster $2$												0.79
cluster $3$												0.80
cluster $4$												0.57
cluster $5$												0.69
cluster $6$												0.75
cluster $7$												0.79
cluster $8$												0.70
cluster 9												0.67
cluster $10$												0.79
cluster $11$												0.61
cluster $12$												0.79
cluster $13$												0.74
cluster 14												0.51
cluster $15$												0.75
cluster 16												0.73
cluster $17$												0.74
cluster $18$												0.49
cluster 19												0.71
cluster 20 $$												0.70
cluster $21$												0.80
cluster $22$												0.68
cluster $23$												0.55
cluster $24$												0.68
cluster $25$												0.77
cluster 26 $$												0.39

Fig. 3. Structures of the clusters for the Letters dataset

because the dataset does not keep the information about the corresponding font for each particular observation.

Cluster		[0, 0.1]	[0.1, 0.2]	[0.2, 0.3]	[0.3, 0.4]	[0.4, 0.5]	[0.5, 0.6]	[0.6, 0.7]	[0.7, 0.8]	[0.8, 0.9]	[0.9, 1]
1	Main class	0	0	0	2.06	25.61	29.33	13.63	9.91	5.78	0
	Other classes	0	0	0	0.21	15.02	16.95	16.3	16.52	21.67	0
2	Main class	0	0	0.81	16.26	24.39	6.5	28.45	10.56	4.06	0
	Other classes	0	0	0	0.36	2.71	10.3	19.16	20.79	19.71	0
3	Main class	0	0	0	0	1.19	15.47	32.14	20.23	17.85	0
	Other classes	0	0	0	0.11	1.15	8.67	14.68	18.26	26.47	0
4	Main class	0	0	0	6.15	19.56	18.11	11.95	10.5	17.02	0
	Other classes	0	0	0	0	0	0	0	4	24	0

5	Main class	0	0	0	0	6.4	21.6	13.6	16.8	18.4	0
	Other classes	0	0	0	0.17	2.89	7.33	12.36	17.81	25.31	0
6	Main class	0	0	0	0	2.55	6.2	14.23	27	25.91	0
	Other classes	0	0	0	0.16	1.63	8.49	11.27	20.09	25.81	0
7	Main class	0	0	11.69	45.61	26.9	12.28	2.33	0.58	0	0
	Other classes	0	0	0	0	6.36	22.61	37.26	16.87	11.78	0
8	Main class	0	0	0	0.86	5.17	19.82	25.86	18.1	16.37	0
	Other classes	0	0	0	1.13	12.3	16.85	16.17	16.62	17.53	0
9	Main class	0	0	25	40.9	13.63	4.54	9.09	6.81	0	0
	Other classes	0	0	21.08	10.88	7.48	8.84	16.32	5.44	16.32	0
10	Main class	0	0	0	1.37	7.9	23.02	27.49	13.74	14.43	0
	Other classes	0	0	0	4.85	12.13	11.4	17.71	23.05	18.68	0
11	Main class	0	0	1.76	16.22	19.57	20.28	14.99	8.81	10.58	0
	Other classes	0	0	0	0	0	0	6.12	14.28	22.44	0
12	Main class	0	0	0	21.56	26.79	26.14	9.8	3.26	4.57	0
	Other classes	0	0	0	2.33	4.67	12.37	20.08	20.22	20.63	0
13	Main class	0	0	0.37	1.5	9.02	22.55	24.06	14.28	12.4	0
	Other classes	0	0	0	0.11	1.95	8.49	13.31	19.51	25.02	0
14	Main class	0	0	0	0	0	0.83	16.66	24.16	26.66	0
	Other classes	0	0	0	0	1.08	7.78	17.65	22.25	25.77	0
15	Main class	0	0.3	23.45	11.72	13.88	13.27	12.65	12.34	8.33	0
	Other classes	0	0	0	0	0	0	0	0	0	0
16	Main class	0	0	0	0	4.91	18.85	17.21	15.57	20.49	0
	Other classes	0	0	0.32	0.75	6.81	14.91	14.27	19.02	21.4	0
17	Main class	0	0	0	0	0	18.85	38.28	25.42	10.28	0
	Other classes	0	0	0	5.75	10	14.24	14.24	14.24	18.18	0
18	Main class	0	0	0	12.66	19.65	16.15	25.32	11.35	6.98	0
	Other classes	0	0	0.27	0.27	7.52	14.76	17.27	19.49	22.84	0
19	Main class	0	0	0	0	4.58	22.32	26.6	20.48	15.59	0
	Other classes	0	0	0	0.2	3.34	6.69	15.06	20.08	23.84	0
20	Main class	0	0	0	0.32	2.96	20.39	22.03	20.72	16.77	0

_	Other classes	0	0	0	0.14	4.41	17.05	18.67	17.79	17.94	0
21	Main class	0	7.1	20.71	32.54	17.75	14.2	5.91	1.77	0	0
	Other classes	0	0	0	0	0	0	0	0	0	0
22	Main class	0	0	0	0	0	2.36	17.32	22.04	31.49	0
	Other classes	0	0	0	0	1.52	9.87	16.68	20.79	26.43	0
23	Main class	0	0	0	1.72	12.06	13.21	22.41	25.28	12.64	0
	Other classes	0	0	0	0.49	2.45	10.49	22.62	30.49	12.95	0
24	Main class	0	0	0	0	3.61	10.04	21.28	21.68	23.69	0
	Other classes	0	0	0	0.24	1.55	6.31	11.88	23.52	27.95	0
25	Main class	0	0	0	0	0.39	3.92	7.45	15.68	43.92	0
	Other classes	0	0	0	0.19	3.16	7.27	13.4	20.3	28.06	0
26	Main class	0	5.35	14.88	20.83	17.85	9.52	10.71	9.52	4.16	0
	Other classes	0	0	0	0	0	0	0	0	0	0

Table 8: Letters: repartition of the classes by cluster layers

Table 9 presents the average depth of each class in each cluster. All the classes are at the same depth inside the clusters. Moreover this depth generally between 0.5 and 0.9 - shows that the points are quite far from the boundary of each cluster. An interesting case is the 14-th cluster: most classes of the dataset are very deep inside the cluster. This can be seen as a strong belonging to the clusters, and it can be concluded that although the cluster is strongly constituted, it does not permit to discriminate between classes. A similar situation appeared for the Australian credit dataset (subsection 2.6).

# 4.3 Summary

- The relationships between classes and clusters are not very explicit
- The points are not as "deep" inside the clusters as they were for the Pendigits dataset. However, they are "deep" enough for most of the clusters and therefore clustering is considered of good quality.
- For most of the clusters the correctly classified points are "deeper" than misclassified points.
- The clusters do not discriminate the classes. A similar situation appeared for the Australian credit dataset. A possible hypothesis is that the clusters grouped the records according to their font recognition rather than letter recognition.

26	ı	ı	0.92	0.96	0.9	0.94	ı	I		ı	0.97	I	0.74	0.78	I	0.91	0.65	0.69	0.89	0.91	ı	0.87	ı	0.8	0.74	0.5	
25		0.83	0.73	ı	0.71	ı.	ı	0.77		ī	I	0.89	0.91	0.77	ı	0.65	0.75	I	ī		ı	ı	ı	0.91	0.83	Т	
24	.84	.82	.76	I	.89	ı	I	.71	I	I	.85	.75	.79	.78	I	.93	I	I	.97	.78	I	.96	ı	.76	.77	ı	
23	.86 (	.87 (	.88 (	ı		.85	ı	.83 (	ı	I	-	.86 (	.83 (	0.8 (	I	-	ı	I	.89 (	.74 (	ī	0.8 (	.83	.86 (	.77 (	ī	
22	-	- 0	0.9		89.	.77 0	ı	.73 0	ı	.(8	I		.85 (	.76	I	.96	ı	I	-	.85 (	ı	.82	0.7	.71 0	.83	ı	
21			.85 (		- 0	.85 0	1	.6 0	.71	.91 (	I	.0	.79 0	.86 0	I	- 0	1	I	77.	.82 0	.38	.73 0	.0 (	.89 0	- 0	1	
0		1	- 0	1		- 0	1	68 (	89 0	- 0	I		- 0	$91 \ 0$	I	1	1	1	.7 0	73 0	- 0	0	97 (	0			
9			6.	93	86	95		- 0.	- 0.		94			83 0.		88	۲.	.6	0	98 0.		66	0.	92	66		
8			85 0.	0.	98 0.	72 0.		- 62		<u>.</u> .	0.		92 -	81 0.		0.	0	0	91 .	97 0.		8 0.	39	95 0.	84 0.		
7 13	- 20	.3	33 0.8	1	8 0.9	9 0.7	1	0.7		0.(	I	4 -	4 0.9	8 0.8	I	5 -	1	5 -	0.9	0.9	1	.0 0.	33 0.6	31 0.9	32 0.8	1	
1	0.5	0.7	1 0.8	1	8 0.	0.	I	2 -	I	I	י פ	0.7	3 0.9	5 0.7	2 -	3 0.7	- 2	8 0.5	, w	- 2	I	7 0.6	0.8	3 0.8	8 0.8	I	
16	'	1	7 0.8	1	8 0.8	1	I	5 0.9	I	I	3 0.9	י מ	8 0.8	3 0.7	0.5	7 0.9	0.7	0.7	0.7	5 0.9	I	0.9	ľ	7 0.7	8 0.6	I	
15	'	ı	1 0.8	I	9 0.8	I	I	2 0.7	I	I	0.8	0.6	0.7	0.8	I	0.9	I	I	-	1 0.9	I	I	T	3 0.7	0.8	I	
14	ı	1	0.7	I	0.86	1	I	.0.6	I	I	I	0.51	0.69	0.8	I	1	ı	I	0.95	0.8	ı	'	'	0.73	0.9	ı	
13	0.6	0.65	0.85	I	0.72	0.92	I	0.79	I	I	I	I	0.92	0.65	I	0.65	0.9	I	ı	0.95	I	0.76	0.85	0.91	0.8	Ţ	
12	ī	0.87	0.78	ı	0.8	0.91	ı	0.72	0.91	I	0.96	0.89	0.88	0.81	I	0.78	ı	I	0.99	0.89	ı	0.74	0.82	0.84	0.84	ī	
11	ı	ı	0.84	ı	0.88	ı	I	0.94	I	I	0.59	I	I	0.84	T	0.95	0.9	I	0.8	0.77	I	0.99	ı	1	0.79	I	
10		ī	0.89	0.89	ī	0.87	ı	0.68	ı	0.68	T	ı	0.99	0.8	I	ı	ı.	I	0.9	0.88	T	0.85	0.71		0.93	ī	
6	ī	ī	0.75	ı	0.77	ı	ı	0.79	ı	ı	0.84	0.7	0.77	0.82	I	0.74	0.95	I	0.8	0.83	ı	ı	ı	0.85	0.83	ı	
$\infty$	ı	ı	0.91	ı	ı	ı	0.61	0.71	0.43	ı	0.88	I	I	0.76	I	I	ı	ı	0.78	0.69	ı	0.91	0.97	ı	0.92	ı	
4	ī	0.91	0.72	ı	0.74	ī	1	0.7	1	I	I	0.71	0.85	.77	I	0.77	0.92	I	ı	0.84	ī	I	I	).83	0.86	ī	
9	ı		0.0	I	.95 (	).93	.41	.67	).29	ı	.93	.81 (	0.7	.81	I	.95 (	1	ı	.79	).73 (	ı	.95	.88	.78	.91	ı	
ъ		ī	.82	.92	.86 (	.79 (	-	.88 (	-	I	-		ı	.74 (	I	.86 (	0.9	.93	.86 (	.83 (	ı	.82 (	.74 (		.72 (	ī	
4	ī	ı	.77 (	-	.79 (	.89 (	.74	.76 (	.45	I	.92	0.6	.77	.83 (	I	.68 (	ı	-	.98 (	.82 (	.98	.77 0	.92 (	.73	.82 (	ı	
en en	.68	.74	.8 0	ı	.74 0	.86 0	- 0	.93 0		ı	.85 0		0	.83 0	I	.72 0	1	ı	0	- 0	- 0	.84 0	0	-	.78 0		
5	0 -	- 0	73 (	29	82 0	88 0		83 0		1	- 0	84	6.	81 0	1	87 0	95	87		93	1	0 6.0	84	93	77 0		
_	63	59	87 0.	- 0.	83 0.	- 0.		- 0.			92	0.	0	88 0.		83 0.	- 0.	- 0.		- 0.		6.	0.	.0	81 0.		
	1 0.	2 0.	3 0.	4.	5 0.	. 9	. 2	∞	6	10	11 0.	12	13	14 0.	15	16 0.	17	18	19	20	21	22 0.	53	24	25 0.	. 56	
					·				ı	· · ·		i		I				· · ·	⊢ ′ ' <b> </b>	••	••	••	••	••	••	••	

Table 9. average depth for each class in each cluster



Fig. 4. The structures of the clusters

• The lowest dissimilarity function value has been reached by the step by step clustering method. The *k*-means method produced a higher value for the dissimilarity function.

## 5 Conclusions

1. The main conclusion to this research is that clusters do not necessarily coincide with classes. Several reasons can cause such results:

- the chosen clustering model does not match the classification structure. (For example, it is possible that some classes in Letters dataset are not point based; at the same time we use the point based model for clustering.)
- the dataset contains a high proportion of noisy records and/or possible mistakes, appeared, for example, on the stage of data collection;
- some characteristics link points more strongly than their belonging to classes.

2. The notion of purity cannot be always used for evaluation of accuracy of clustering methods. If the purity is high enough, we can conclude that the chosen clustering method is efficient for the dataset under consideration. However, we can not make any conclusion regarding the efficiency of clustering methods, if purity is low. In such a case the classes and clusters do not coincide. There can be different reasons for this: either the applied clustering technique works not very well, or there might be some other "hidden" characteristics, which link the records together and which are independent from the classes.

Some other approaches (different from comparing classes and clusters) can be used for evaluating the quality of clustering:

- the value of a cost function (called *dissimilarity function*, and in our case *cluster function*);
- the structure of the clusters (how "deep" the points are inside the clusters).

## References

- A. M. Bagirov, Derivative-free methods for unconstrained nonsmooth optimization and its numerical analysis, Investigacao Operacional, Vol. 19, 1999, pp 75-93.
- [2] A. M. Bagirov, Numerical methods for minimizing quasidifferentiable functions: a survey and comparison, In: V.F. Demyanov and A.M. Rubinov (eds.), Quasidifferentiability and Related Topics, Kluwer Academics Publisher, pp 33-71, 2000.
- [3] A.M. Bagirov, A.M. Rubinov and J. Yearwood, A global optimization approach to classification, *Optimization and Engineering*, 3, 2002, 129-155.
- [4] Adil M. Bagirov and Julien Ugon, An algorithm for minimizing clustering functions, Optimization, accepted for publication.

- [5] A. M. Bagirov and A. M. Rubinov, Global minimization of increasing positively homogeneous function over unit simplex. Annals of Operations Research, Vol. 98, 2000, pp 171-187.
- [6] A. Bagirov, A. Rubinov, N. Soukhoroukova, J. Yearwood, Unsupervised and Supervised Data Classification Via Nonsmooth and Global Optimization, Top, Vol. 11, Number 1, 1-93, Sociedad de Estadistica Operativa, Madrid, Spain, June 2003.
- [7] I.S. Dhillon, J. Fan and Y. Guan, Efficient clustering of very large document collections, In: *Data Mining for Scientific and Engineering Applications*, Kluwer Academic Publishers, Dordrecht, 2001.
- [8] A. K. Jain, M. N. Murty and p. J. Flynn Data Clustering: A Review ACM Computing Surveys, Vol. 31, No 3, September 1999.
- J. B. MacQueen, Some Methods for Classification and Analysis of Multivariate observations. In: L. M. LeCam and J. Neyman (eds.), Proceedings of the Firth Berkeley Symposium on Mathematical Statistics and Probability, Vol. 1. Berkeley: University of California Press, 1967.
- [10] D. Michie, D. J. Spiegehalter and C. C. Taylor (eds.), Machine learning, neural and statistical classification. Ellis horwood series in artificial intelligence, London, 1994.
- [11] B. Mirkin, Mathematical Classification and Clustering, Kluwer Academic Publishers, 1996.
- [12] A. Rubinov, Abstract Convexity and Global Optimization, Kluwer Academic Publishers, 2000, ISBN 0-7923-6323-X.
- [13] A. Rubinov, N. Soukhoroukova, J. Yearwood, Clustering for studying the structure and quality of datasets, L. Caccetta, V. Rehbock, Industrial Optimisation, Vol 1, Proceedings of Symposium in Industrial Optimisation, Western Australian Centre of Excellence in Industrial Optimisation (WACEIO), Curtin University, Perth, Australia, 2003.
- [14] A. Rubinov, N. Soukhoroukova, J. Ugon, *Minimization of the sum of minima of convex functions and its application to clustering*, to appear in Continuous optimization: current trends and applications, Springer.
- [15] http://www.ciao-go.com.au/index.php