



Scaling-up multiobjective evolutionary clustering algorithms using stratification

Alvaro Garcia-Piquer^a, Jaume Bacardit^b, Albert Fornells^{c,**}, Elisabet Golobardes^d

^aInstitute of Space Sciences (IEEC - CSIC), Campus UAB, Faculty of Science, C5 Tower - 2nd Floor, 08193 Bellaterra, Spain

^bSchool of Computing Science, Newcastle University, Claremont Tower, Newcastle, NE1 7RU, UK

^cResearch Group in Hospitality, Tourism and Mobilities, HTSI, Universitat Ramon Llull. Av. Marquès de Mulhacén 40-42, 08034 Barcelona, Spain

^dGR-SETAD, La Salle, Universitat Ramon Llull. Av. Quatre Camins 30, 08022 Barcelona, Spain

ABSTRACT

Multiobjective evolutionary clustering algorithms are based on the optimization of several objective functions that guide the search following a cycle based on evolutionary algorithms. Their capabilities allow them to find better solutions than with conventional clustering algorithms when more than one criterion is necessary to obtain understandable patterns from the data. However, these kind of techniques are expensive in terms of computational time and memory usage, and specific strategies are required to ensure their successful scalability when facing large-scale data sets. This work proposes the application of a data subset approach for scaling-up multiobjective clustering algorithms and it also analyzes the impact of three stratification methods. The experiments show that the use of the proposed data subset approach improves the performance of multiobjective evolutionary clustering algorithms without considerably penalizing the accuracy of the final clustering solution.

© 2016 Elsevier Ltd. All rights reserved.

1. Introduction

Multiobjective clustering (MC) algorithms (Handl and Knowles, 2007) tackle the challenge of optimizing several criteria that cannot be combined in a single objective function by defining an objective function for each criterion and by optimizing them trying to obtain a trade-off between all the objectives. There are different strategies for multiobjective optimization such as Simulated Annealing (Saha and Bandyopadhyay, 2010) and Ant Colony Optimization (Iredi et al., 2000), but Multiobjective Evolutionary Algorithms (MOEAs) (Coello, 1999) have become one of the most capable strategies to solve this kind of problem (Zitzler et al., 2000) since they (1) work with a collection of solutions with different trade-offs among objectives, which are improved until a Pareto set with optimal trade-offs is obtained; (2) can be easily adapted to the type of data of the studied domain, due to the flexible knowledge representation used; and (3) are able to optimize different objectives without assuming any underlying structure of the objective functions. Therefore, MOEAs offer outstanding search capabilities but their performance can be compromised in large databases

due to their high computational and memory usage requirements (Freitas, 2002). This is an important issue considering the needs when analyzing large data sets in a reasonable computational time and memory usage without considerably penalizing their accuracy (Kargupta et al., 2009). Evolutionary Algorithms (EAs) are based on the principles of evolution and natural selection, applying an iterative process where a collection of initial solutions (individuals) are evolved through pseudo-random recombination until obtaining an optimal solution. In clustering, an individual is a possible group of data. From the whole process, the evaluation is the most time-consuming step because each individual has to be assessed with respect to all the elements according to all the objective functions defined.

This work proposes to scale-up MOEAs when they are applied to large data using an approach based on data subsets techniques. Specifically, the approach splits the data set into several strata so the EA only uses a stratum for evaluating the individuals in each generation following a Round Robin policy to avoid bias problems. Thus, computational and memory costs associated to the evaluation of the population are drastically reduced and its application does not need to modify the algorithm structure. An ideal stratification strategy is to map the initial data set into disjoint strata of equal size and with equal class distribution and where the number of strata is defined by the user. How-

**Corresponding author: Tel.: +34-932522890
e-mail: albert.fornells@htsi.ur1.edu (Albert Fornells)

ever, clustering problems are unsupervised and classes cannot be used to split the instances into representative strata because they are unknown (Cano et al., 2008). Therefore, the definition of the size of subsets and the selection of their elements are not trivial and they influence the performance of the algorithm if they are not sufficiently representative. For this reason, two unsupervised and one supervised strata generation strategies are presented and analyzed: (1) random strata, (2) strata according to clusters distribution using a fast and approximate clustering algorithm and (3) strata based on classes. Moreover, strategies are integrated in a MC algorithm based on MOEAs called CAOS (Clustering Algorithm based on multiObjective Strategies) (Garcia-Piquer et al., 2012, 2013) and they are tested with several strata size using artificial and real world problems. Finally, results are compared with each other and with regard to MOCK, which is one of the most well-known MC algorithms based on MOEAs (Handl and Knowles, 2007). The results show improvement in computational time while accuracy is not substantially penalized when stratification approaches are applied. Furthermore, the three strategies for building the strata are equivalent so the proposed data subset approach can be used in clustering problems because it does not need a stratification method based on classes. Finally, the results also show that the proposed approach is significantly better than MOCK in terms of computational time and accuracy.

The paper is organized as follows. Section 2 summarizes the related work on data subsets applied to clustering. Sections 3 and 4 describe CAOS and the stratification strategies. Section 5 describes the experimentation and discusses the results. Finally, Section 6 ends with conclusions and further work.

2. Related Work

Two of the most used strategies for scaling-up EAs are Parallel EA (Cantu-Paz, 2000) and data subsets (Cano et al., 2008; Derrac et al., 2010). The first strategy distributes the computational cost of the evaluation step by parallelizing the evaluation of individuals so it is necessary to adapt or redefine the algorithm in order to be able to parallelize it in an environment with several processors. Moreover, the parallelization may imply an additional communication cost that could decrease the performance achieved with the distribution of compute. On the other hand, the second strategy uses a data subset from the original data set to evaluate the individuals so fewer resources are required and there is no need to modify the algorithm structure. In contrast, the data sets definition is not trivial.

There are two main ways to work with data subsets: using only one of the built data subsets, or using alternatively all the data subsets. The algorithm CLARA (Clustering LARge Applications) (Kaufman and Rousseeuw, 1990), one of the most representative algorithms for clustering large data sets, works using the first approach. This algorithm is based on selecting randomly a sample from the entire data set and, subsequently, it finds k medoids of the sample using only the built sample. After this, all the instances of the entire data set are assigned to the most similar medoid. The execution of the entire process is repeated five times, and the solution with less dissimilarity

is returned as the solution. Following this idea, other methods consist of randomly extracting several samples from the entire data set and applying the same clustering algorithm to each one of the samples, thereby obtaining several clustering results. After this, all the obtained results are merged in a single clustering solution. Hore et al. (Hore et al., 2009) proposed using k -means or fuzzy k -means algorithms with large data. The idea is to obtain a set of jointed or disjointed samples and apply one of the two algorithms to each sample to obtain several clustering results. The last step consists of doing a consensus between each clustering result to obtain a final clustering solution as in ensemble clustering. The drawback of using only one sample to obtain the clustering results is that it is necessary to execute the algorithm several times or apply it to different data subsets in order to avoid the bias of using only one sample. Moreover, only a part of the entire data set is used. Thus, the approaches based on using all the data subsets can be useful to obtain the clustering results in a single execution.

ILAS (Incremental Learning by Alternating Strata) (Bacardit, 2004) is a technique based on Evolutionary Algorithms for classification problems based on dividing the training set into several strata based on using a different stratum in each iteration of the evolutionary algorithm using a round-robin policy. Thus, the individuals are evaluated with all the strata, avoiding any bias of the data and increasing the generalization of the individual. The strategy followed in this paper is based on the ILAS algorithm but applied to MC problems.

3. CAOS

CAOS (Garcia-Piquer et al., 2013) is a multiobjective evolutionary algorithm system to solve clustering problems based on adapting the multiobjective optimization algorithm PESA-II (Corne et al., 2001) due to its competitiveness with respect to the state-of-the-art clustering methods and its ability to evolve accurate clusterings from domains with complex structures (Handl and Knowles, 2007). It evolves a set of mutually non-dominated clustering solutions (called Pareto set) that correspond to different tradeoffs between objectives. A solution S is non-dominated when there is not any solution better than S in all the objectives. Otherwise, the solution is dominated.

Algorithm 1 summarizes the main elements of PESA-II. It evolves an external population (EP) of individuals through a number of generation where individuals are selected, crossed and mutated following the typical evolutionary cycle. Individuals are represented with real numbers that represent the coordinates (attributes) of the cluster prototype using a centroid-based representation (Hruschka et al., 2009). More specifically, each individual consists of $n \cdot t$ genes $\{x_{11}, \dots, x_{1t}, \dots, x_{n1}, \dots, x_{nt}\}$, where n is the number of clusters of the individual, t is the number of the attributes of the data set, and x_{ij} is the value of the attribute j of the cluster centroid i . The genotypic representation is transformed into the phenotypic representation by assigning each instance to the cluster with the nearest centroid to it. In addition to EP , it also maintains an *internal population* (IP) to separate the exploration from the storage of the best solutions. That is, IP is used to explore new promising solutions and EP is employed to store a large and diverse set of non

```

1 Let  $EP$  and  $IP$  be an external and an internal population respectively.
  They store a maximum of  $N_{EP}$  and  $N_{IP}$  individuals, where ( $N_{IP} < N_{EP}$ )
2 Init.  $IP$  with  $N_{IP}$  individuals stochastically created
3 Init. the  $EP$  individuals with non-dominated clustering results from  $IP$ 
4 Evaluate all the individuals from  $EP$  according to the objectives
5 foreach Generation do
6   Select  $N_{IP}$  individuals from  $EP$  to generate a new  $IP$ 
7   while ( $|IP| \neq 0$ ) do
8     Select and remove two individuals from  $IP$ 
9     Cross and mutate them to obtain 2 new ind.:  $I_{New_1}$  and  $I_{New_2}$ 
10    foreach  $I_{New_i}$  do
11      Evaluate the  $I_{New_i}$  fitness according to the objectives
12      if  $I_{New_i}$  dominates any individual from  $EP$  then
13        Remove the dominated individuals by  $I_{New_i}$  from  $EP$ 
14        Add  $I_{New_i}$  into  $EP$ ;
15      else if  $I_{New_i}$  is not-dominated and  $I_{New_i}$  not-dominates any
16      individual then
17        if  $EP$  is full then
18          Remove an ind. from the most crowded niche
19          Add  $I_{New_i}$  into  $EP$ 
19 Select a individual from  $EP$  as a solution

```

Algorithm 1: Scheme of PESA-II algorithm.

dominated solutions found so far. Moreover, EP is organized in N_{niches} different niches through the placement of an hyper-grid in the objective space splitting it in hyper-rectangles, where each of them is considered as a separate niche. Therefore, solutions with similar objectives will be placed in the same niche. The replacement process uses the niching mechanism to make pressure toward balancing the allocation of solutions in different niches thus encouraging solutions to cover all the objective space. Concretely, the system creates IP with N_{IP} individuals stochastically initialized. All the non-dominated solutions of IP are used to build EP .

The final step is to select a solution from the Pareto set (composed by all non-dominated solutions) when the evolutionary process ends. This point is not trivial because there is not any single individual which is the best in all the objectives and for this reason clustering validation techniques (Halkidi et al., 2001) are required for selecting the best one. CAOS integrates some supervised and unsupervised techniques to score a solution. The supervised approach follows the idea that similar elements from the same class should be in the same cluster. In contrast, the unsupervised approaches retrieve the best solution from the Pareto set according to the quality of the clusters based on the compacting and separation between them, such as Davies-Bouldin index, Dunn index or Silhouette index, among others. Experimentation of Section 5 uses a supervised validation technique based on classes called Adjusted Rand Index (Yeung and Ruzzo, 2001) for being able to compare the performance of several strategies between them. It is important to highlight that class information is not used in the building process of the clusters, it is only used at the end of the process for comparing the obtained clustering solution between them. The Adjusted Rand Index compares a clustering solution with respect to an ideal partition of the data set by counting the number of pairwise co-assignments of instances between them and introducing a statistically induced normalization in order to yield

values close to 0 for random partitions. In contrast, a value of 1 means that all the clusters correspond to the structure defined in the ideal partition. In our case, the ideal partition corresponds to the classes of the instances in the original data set.

Finally, CAOS allows the customization of the objectives, genetic operators among other features. *Deviation* and *Connectivity* are used as objective functions (Handl and Knowles, 2007) because they indicate how nearby are the elements of each cluster (intra-cluster variance) and how separated are the clusters between them (inter-cluster variance), respectively. Both objectives have to be minimized because the desired solution should contain compact clusters with examples that are close in the feature space. A one-point crossover operator (Goldberg, 2002) is used to generate two offspring from pairs of parents. One point is selected for each parent and the parts are interchanged between them, taking into account that they have to cut the individuals at the same attribute but not necessarily at the same cluster. This is an easy crossover strategy according to the size of each individual can be different. A cluster-oriented mutation operator (Hruschka et al., 2009) is used to promote the right search. It defines three different types of mutations and all of them have the same probability to be applied: to merge two clusters, to split a cluster, and to move the centroid of a cluster.

4. Data Subset Strategies

Reducing the amount of data used by an algorithm is a smart approach to reduce the computational cost of evolutionary-based machine learning techniques and it could also improve the accuracy of the system Bacardit (2004). In this sense, we want to scale-up a MC algorithm based on EAs by dividing a data set in several stratified subsets and using them alternatively during the algorithm process in order to avoid bias. Next points detail the analyzed approaches to build the strata and how to use it in MC. Finally, the impact of using these strategies in terms of computational cost and memory usage is described.

4.1. Creation of Strata

Data subset strategies map the initial data set into disjoint data subsets (strata) of equal size and with equal class distribution (Bacardit, 2004; Cano et al., 2006), where the number of strata is selected by the user (see Algorithm 2). Classes in clustering problems are usually unknown, so the stratification based on classes is not always available. To avoid this lack, two unsupervised approaches to divide the data set are proposed:

- **Random Strata.** It randomly assigns the instances to each one of the strata as Algorithm 3 shows.
- **Strata based on Clusters.** It uses a fast and approximate clustering technique to create a partition of the original data set. Next, the data set is stratified according to the obtained clusters, that is, it assigns the instances to each stratum respecting in it the same cluster distribution of the instances than in the clustered original data set. The process is described in Algorithm 4. The clusters are found with the Subtractive Clustering algorithm Chiu (1994) applied to the original data set, which is an efficient and non-iterative method for estimating cluster centers. It is usually

used to determine the number of clusters and their initial values for initializing iterative optimization-based clustering algorithms. The limitation of this strategy is their computational cost $O(m^2)$, where m is the number of instances of the original data set, because with very large data sets can be expensive in computational terms. Nevertheless, in terms of spatial cost it only needs the data set information and a list with the prototypes of each cluster.

Finally, these approaches require the definition of the number of strata which will influence the performance. As the number of strata increases the computational time decreases but pattern extraction becomes more complex due to the lack of information. It is important to highlight that the idea of these strategies is to obtain data with similar distribution in each stratum, and this only can be possible if the size of each one is not very small.

4.2. Evolution Based on Strata

The idea is to use a different stratum in each iteration of the evolutionary algorithm in order to avoid the bias produced when only one stratum is used. Thus, the final individuals can generalize more than using only one of the strata and there is no need to modify the main process. More precisely, the generation of strata is done before line 1 of the Algorithm 1, and the change of stratum is done between lines 6 and 7. It is important to highlight that these strategies can be applied to CAOS due to the fact that individuals are represented by the prototypes of the clusters and then the individuals are independent of the instances. Thus, the algorithm can work with different instances of the data set in each iteration.

```

1 Let numStrata be the number of strata to generate
2 Let Strata be a vector of size numStrata where each position is initially
  an empty list of instances
3 Let I be a vector of size numClasses where each position stores a list of
  the instances of the same class
4 stratum = 0
5 class = 0
6 while (class < numClasses) do
7   while (|I[class]| != 0) do
8     Select randomly an instance i from I[class]
9     Add i to Strata[stratum]
10    Erase i from I[class]
11    stratum = (stratum + 1) mod numStrata
12   class = class + 1
13 return Strata

```

Algorithm 2: Strata generation based on classes.

```

1 Let numStrata be the number of strata to generate
2 Let Strata be a vector of size numStrata where each position is initially
  an empty list of instances
3 Let I be a list of all the instances of the data set
4 stratum = 0
5 while (|I| != 0) do
6   Select randomly an instance i from I
7   Add i to Strata[stratum]
8   Erase i from I
9   stratum = (stratum + 1) mod numStrata
10 return Strata

```

Algorithm 3: Strata generation based on random instances selection.

4.3. Computational Performance Models

The aim of the data subsets strategies is to reduce the computational time and memory usage without considerably penalizing the accuracy. Next, this subsection analyzes the improvement in the performance of a MC algorithm based on EAs using CAOS from a theoretical perspective and section 5 will analyze the improvement using different real and artificial data sets.

The CAOS process can be divided in two main blocks: the initialization process and the clustering process. The initialization process focuses on precalculating the distances between the instances and the nearest neighbors to speed-up the clustering process avoiding the repetition of calculations. In contrast, the clustering process is referred to the evolutionary cycle that obtains the Pareto set of solutions. CAOS algorithm has the same initialization and clustering cost independently of the data subset strategy used (based on classes, random or based on clusters) but it depends on the number of strata used. However, the time of both processes is extremely lower in comparison with the time spent when the complete data set is used as Table 1 describes. It should be emphasized that both times are reduced when the size of the stratum is decreased, that is, when the number of strata increases. Nevertheless, the use of strata requires an additional cost for building them. According to this, the strategies based on random instances selection and based on classes need only one data scan to build the strata and their cost is $O(m)$, where m is the number of instances of the complete data set. In contrast, the strategy based on approximative clusters has a higher cost due to the cost related to the subtractive clustering technique ($O(m^2)$).

As it has been explained above, CAOS precalculate the distances between all the instances of the data set and the nearest neighbors of each instance to speed-up the clustering process. Thus, the memory usage would be extremely high if the complete data set is analyzed when a large data set is used. Applying any of the three strategies of data subset construction, the memory usage is considerably reduced as Table 2 shows. Even the computational time and the memory usage of the MC algorithm is considerably reduced, the accuracy of the method can be penalized due to the fact that less information is used to obtain the clustering solutions.

```

1 Obtaining the instances clustered in numClusters clusters by applying the
  Subtractive Clustering algorithm to the complete data set
2 Let numStrata be the number of strata to generate
3 Let Strata be a vector of size numStrata where each position is initially
  an empty list of instances
4 Let I be a vector of size numClusters where each position stores a list of
  the instances assigned to the same cluster
5 stratum = 0
6 cluster = 0
7 while (cluster < numClusters) do
8   while (|I[cluster]| != 0) do
9     Select randomly an instance i from I[cluster]
10    Add i to Strata[stratum]
11    Erase i from I[cluster]
12    stratum = (stratum + 1) mod numStrata
13   cluster = cluster + 1
14 return Strata

```

Algorithm 4: Strata generation based on approximative clusters.

Table 1. Computational cost of CAOS applied to the complete data set and to data subsets (CAOS_{CD} and CAOS_{DS} respectively) broken down in initialization cost and clustering cost. Where g is the number of generations, $|P|$ is the internal population size, m and t are the number of instances and attributes of the data set respectively, \bar{n}_{cd} is the average of the number of clusters of the individuals (the minimum number of clusters is 1 and the maximum m), \bar{n}_{ds} is the average of the number of clusters of the individuals (the minimum number of clusters is 1 and the maximum $\frac{p}{numStrata}$), ℓ is the percentage of the nearest elements taken into account, and $numStrata$ is the number of strata generated.

Algorithm	Initialization cost	Clustering cost
CAOS _{CD}	$O(m^3 \cdot \ell)$	$O(g \cdot P \cdot m \cdot \bar{n}_{cd} \cdot t)$
CAOS _{DS}	$O(numStrata \cdot (\frac{m}{numStrata})^3 \cdot \ell)$	$O(g \cdot P \cdot \frac{m}{numStrata} \cdot \bar{n}_{ds} \cdot t)$

Table 2. Memory usage of CAOS applied to the complete data set and to data subsets (CAOS_{CD} and CAOS_{DS} respectively) to store the nearest neighbors. Where m is the number of instances of the data set, ℓ is the percentage of instances considered neighbors, $numStrata$ is the number of strata generated and $sizeof(data\ type)$ is the size in bytes of the data type.

Algorithm	Storage of distances	Storage of nearest neighbors
CAOS _{CD}	$m^2 \cdot sizeof(float)$	$(\ell \cdot m)^2 \cdot sizeof(integer)$
CAOS _{DS}	$(\frac{m}{numStrata})^2 \cdot sizeof(float)$	$numStrata \cdot (\frac{\ell m}{numStrata})^2 \cdot sizeof(integer)$

5. Experiments, Results, and Discussion

This section evaluates the performance improvement in CAOS using the proposed data subset approach with the three stratification strategies described in section 4. The performance is considered in terms of accuracy using the Adjusted Rand Index value of the solution returned by CAOS and the computational time required to find it.

The organization of this section is as follows. First, the collection of 105 artificial and real-world data sets and the experimental methodology are described. Next, CAOS results and a comparison between the proposed data subset approach in CAOS and MOCK (Handl and Knowles, 2007) is presented in order to emphasize the performance improvement of our approach. Finally, the performance of the most suitable strategy is analyzed in data sets of medium and large size.

5.1. Test Bed

The experimentation uses different typologies of artificial and real-world problems. Concretely, 75 artificial data sets were created according to different number of instances (from 800 to 24000), attributes (from 2 to 100) and classes (from 2 to 30). They were built adapting the tool used in (Handl and Knowles, 2007) where three parameters are used to create the data sets: the number of attributes, the number of classes related to the number of instances, and the separation between the classes. Each class has a data distribution for each attribute, which can only have numerical values. The distribution can be a normal or uniform distribution, and it is randomly selected to model each attribute. Also, the separation between classes were modeled, obtaining 25 data sets with well-separated classes, other 25 data sets with nearer classes, and the last 25 with overlapped classes. On the other hand, other 30 real-world problems were selected according to different number of instances (from

150 to 58000), attributes (from 2 to 60) and classes (from 2 to 26) from the UCI repository (Frank and Asuncion, 2010). The characteristics of each data set are detailed in Table 5 and Table 6 in supplementary material. It must be emphasized that the class assigned to each instance of the data sets is known in order to apply the stratification strategy based on classes and to evaluate the accuracy of the clustering results.

5.2. Experimental Methodology

The performance of the three approaches based on data subsets (CAOS_{DS}) were compared with respect to the approach that uses the complete data set (CAOS_{CD}) in terms of accuracy and computational time. The accuracy is compared using the Adjusted Rand Index, which is based on the initial classes of the data set, where 1 is the best accuracy (all the clusters correspond to the original classes) and 0 the worst. The computational time represents the sum of the precalculation time and clustering time. The first one includes the time needed to build the data subsets and to precalculate the distance and nearest neighbors structures necessary to the clustering process. The second one is referred to the time needed to do the evolutionary process that obtains the Pareto set of solutions. Finally, each CAOS_{DS} strategy is executed dividing the original data set in 2, 3, 4, 10, 15, 20 and 25 data subsets which means the 50%, 34%, 25%, 20%, 10%, 7% and 4% of the instances of the original data sets are considered in each data subset respectively.

Each CAOS configuration was run with 20 different seeds and with the following parameters (see Section 3 for notation details): ℓ is 5% of the number of instances used, the maximum size of the initial population is 100, N_{EP} is 1000, N_{IP} is 50, N_{niches} is 5, the number of generations is 400, P_c is 0.7 and P_μ is $1/m$. The minimum and maximum number of clusters for the initial individuals is 2% and 20% of m respectively.

On the other hand, the recommendations pointed out by Demšar (Demšar, 2006) were followed to perform the statistical analysis of the accuracy of the algorithms, which was based on the use of nonparametric tests. More specifically, we followed the process given in (García and Herrera, 2008) to compare them using the software freely provided by the authors. First, the Friedman's test (Friedman, 1940) with $\alpha = 0.05$ was applied to contrast the null hypothesis that all the learning algorithms obtained the same results on average. Then, if the Friedman's test rejected the null hypothesis, pair-wise comparisons were performed by means of the Holm's step-down procedure (Holm, 1979). Following this procedure, we could distinguish pairs of learners that performed significantly differently.

5.3. Discussion of CAOS Results

The accuracy of CAOS_{CD} and the three CAOS_{DS} strategies were empirically tested with the presented 105 data sets using the Holm's test (see Table 7 in the supplementary material). Analyzing the results it can be observed that independently of the size of the strata, all the CAOS_{DS} strategies are not significantly different between them. However, all of them are significantly different to CAOS_{CD} due to the fact that they are using less information in the clustering process. For this reason, we repeat the same Holm's test but separating data sets in artificial

Table 3. Comparison of the accuracy of the algorithms in the (a) artificial and (b) real data sets using the post-hoc Holm’s procedure with $\alpha = 0.05$. The algorithms compared are CAOS using the complete data set CAOS_{CD} and the three CAOS_{DS} strategies to generate data subsets: based on classes (CAOS_{DS} Classes), random (CAOS_{DS} Random) and based on clusters (CAOS_{DS} Clusters). The results are showed for 34%, 10% and 4% of information used from the complete data set. The symbols \oplus and \ominus show that the method in the row obtained results that were significantly higher/lower than those obtained with the method in the column. Similarly, the symbols + and – denote a non-significant higher/lower results. The Iman and Davenport statistic is calculated according to F-distribution with 3 and 222 degrees of freedom.

(a)							
% Instances	Strategies	CAOS _{CD}	CAOS _{DS} Classes	CAOS _{DS} Random	CAOS _{DS} Cluster	Avg. Rank	p-value
34%	CAOS _{CD}		\oplus	\oplus	\oplus	1.871	1.705E – 5
	CAOS _{DS} Classes	\ominus		+	+	2.662	
	CAOS _{DS} Random	\ominus	–		+	2.689	
	CAOS _{DS} Clusters	\ominus	–	–		2.777	
10%	CAOS _{CD}		\oplus	\oplus	\oplus	1.730	2.634E – 8
	CAOS _{DS} Classes	\ominus		+	+	2.716	
	CAOS _{DS} Random	\ominus	–		–	2.817	
	CAOS _{DS} Clusters	\ominus	–	+		2.736	
4%	CAOS _{CD}		\oplus	\oplus	\oplus	1.676	4.689E – 10
	CAOS _{DS} Classes	\ominus		+	+	2.662	
	CAOS _{DS} Random	\ominus	–		–	2.946	
	CAOS _{DS} Clusters	\ominus	–	+		2.716	

(b)							
% Instances	Strategies	CAOS _{CD}	CAOS _{DS} Classes	CAOS _{DS} Random	CAOS _{DS} Cluster	Avg. Rank	p-value
34%	CAOS _{CD}		\oplus	+	+	1.816	0.005
	CAOS _{DS} Classes	\ominus		–	–	2.917	
	CAOS _{DS} Random	–	+		+	2.650	
	CAOS _{DS} Clusters	–	–	+		2.617	
10%	CAOS _{CD}		\oplus	\oplus	\oplus	1.533	4.853E – 6
	CAOS _{DS} Classes	\ominus		–	–	2.966	
	CAOS _{DS} Random	\ominus	+		–	2.933	
	CAOS _{DS} Clusters	\ominus	+	+		2.567	
4%	CAOS _{CD}		\oplus	\oplus	\oplus	1.400	8.677E – 8
	CAOS _{DS} Classes	\ominus		+	+	2.667	
	CAOS _{DS} Random	\ominus	–		–	3.000	
	CAOS _{DS} Clusters	\ominus	–	+		2.933	

and real-world datasets in order to analyze in more detail the performance of the different strategies as Table 3(a) and 3(b) illustrates using 34%, 10% and 4% of instances for the artificial and real data sets respectively (see Tables 8(a) and 8(b) from the supplementary material for the complete experimentation). The results with respect to the artificial data sets shows that the three CAOS_{DS} strategies are significantly worse in terms of accuracy regarding CAOS_{CD} independently of the number of instances considered. Nevertheless, the three CAOS_{DS} are not significantly different in terms of accuracy between them. From a quantitative point of view, Figure 1(a) shows the average of accuracy difference among the three CAOS_{DS} strategies and CAOS_{CD}. Globally, the strategies based on classes and on clusters follow a similar pattern, and the accuracy is not con-

siderably decreased until less than the 20% of the instances are used. In contrast, the random selection approach seems to underperform the other two due to the fact that the classes structure is complex in some data sets and the random strategy is not able to build representative strata. Figure 2(c) shows the average speedup for each data set of the three CAOS_{DS} strategies. The speedup of each CAOS_{DS} strategy applied to a specific data set corresponds to the time needed to apply CAOS by CAOS_{CD} divided by the time needed by the corresponding CAOS_{DS} strategy. Thus, a high speedup is desired (e.g., a CAOS_{DS} strategy with speedup of 3 indicates that it is three times faster than the CAOS_{CD} approach). This speedup is divided in two partial speedups (1) regarding the precalculation time needed to build the data subsets and to precalculate the distance and nearest

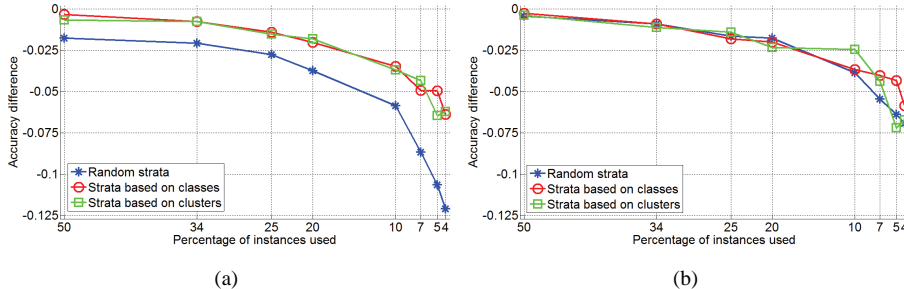


Fig. 1. Accuracy difference of the three CAOS_{DS} strategies regarding CAOS_{CD}. (a) Artificial data sets and (b) real-world data sets.

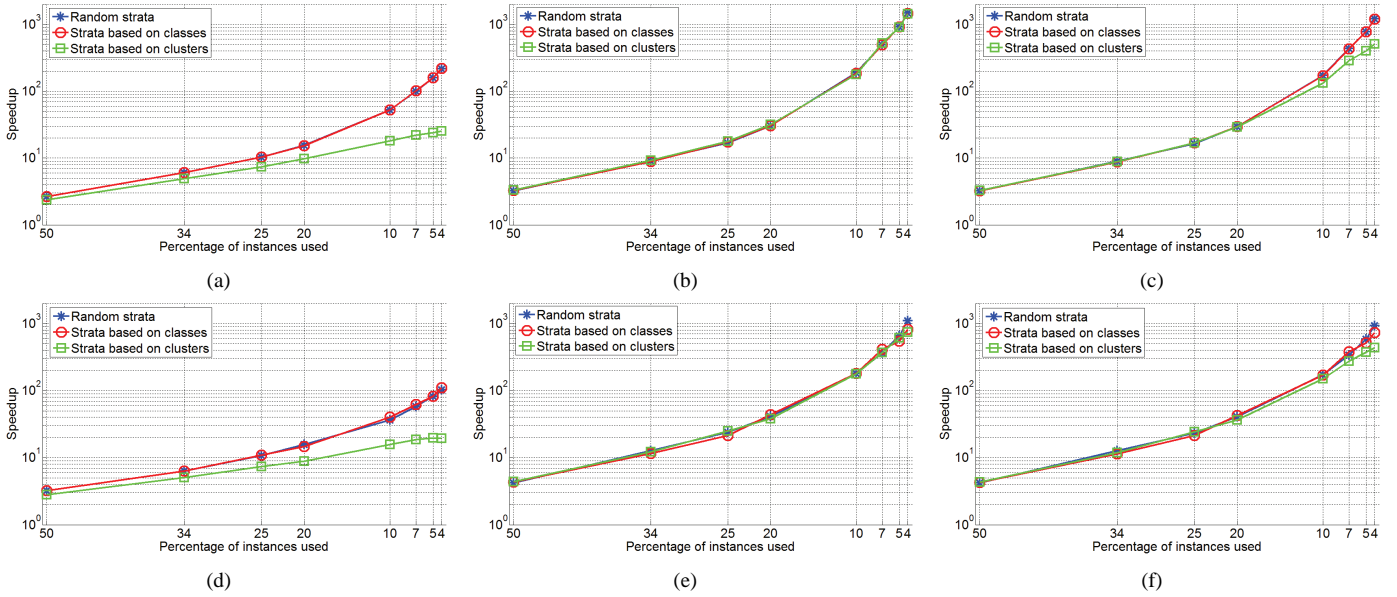


Fig. 2. Speedup of the three CAOS_{DS} strategies regarding CAOS_{CD} in artificial (a,b,c) and real-world data sets (d,e,f). Figures (a,d) are referred to the speedup of the precalculation time. Figures (b,e) show the speedup of the clustering time to do the evolutionary process that obtains the Pareto set of solutions. Figures (c,f) are related to the speedup of the overall time taking into account both times.

neighbors structures (see Figure 2(a)), and (2) regarding to the clustering time (see Figure 2(b)). The precalculation speedup is the same in the random strata than in the strata based on classes strategies. However, the strategy based on clusters has a lower precalculation speedup because it needs to roughly cluster the instances before building the data subsets. On the other hand, the speedup of the clustering time is the same in the three strategies since the strata method does not affect to the clustering process. Analyzing the overall speedup according to the percentage of instances used from the complete data set, it can be observed that using a 50% of the instances, that is the lowest improvement, the speedup is 3, so CAOS_{DS} strategies are three times faster than CAOS_{CD}. Moreover, using a 4% of the instances of the complete data set the speedup is 500 for the strategy based on clusters and 1200 for the other two strategies. Also it can be observed that there are not speedup differences between the strategies until more than a 10% of the instances are used. These analysis showed that CAOS_{DS} strategies applied to the proposed artificial data sets are faster than CAOS_{CD} and do not considerably penalize the accuracy in quantitative terms. Nevertheless, CAOS_{DS} is worst and significantly different, in statistical terms, regarding CAOS_{CD}.

With respect to the results of the Holm's test applied to the real-world data sets, the Friedman's test cannot reject the null hypothesis that all the strategies obtain the same results on average when in each data subset is considered a 50% of the instances (2 data subsets). Thus, the three CAOS_{DS} strategies cannot be considered different than CAOS_{CD} in terms of accuracy. When the 34% and the 25% of the instances are considered, the CAOS_{DS} strategy based on random strata and on clusters is not significantly different in terms of accuracy regarding CAOS_{CD}. Nevertheless, the accuracy of the other strategy is worst and significantly different than CAOS_{CD}. When it is used less than the 25% of the instances of the complete data set in each data

subset, the three CAOS_{DS} strategies are significantly different in terms of accuracy regarding CAOS_{CD}, but they are not significantly different between them. Figure 1(b) shows that there is virtually no accuracy difference between CAOS_{DS} strategies regarding CAOS_{CD}, because the data sets used do not have shapes as complex than the artificial ones used. Figure 2(f) shows that, in terms of speedup, the behavior in the used real-world data sets is similar than in artificial data sets. The maximum speedup is lower because some of the real-world data sets are small and the speedup of using less than a 10% of the instances is not as high than in larger data sets. Also, it can be observed than in the used real-world data sets if it is used a 50% of the instances, it is four times faster than CAOS_{CD} and it obtains the same clustering results (see Table 11 and Table 12 in the sup. material with the results of the execution time). It must be emphasized that the speedup obtained applying this kind of techniques is very high and, consequently, the computational performance of the system is considerably improved. Moreover, assuming that the best strata generation is based on the original classes, the results show that the other two strategies to build the strata are not significantly different in terms of accuracy independently of the kind of data sets tested. In terms of accuracy, the random and the cluster based strategies are as useful as the strategy based on classes but without the requirement of having the original class of each instance. In terms of computational time, the random and classes based strategies have similar speedup regarding CAOS_{CD}. Nevertheless, the cluster based strategy has a lower speedup. According to these observations, it seems that the most suitable strategy to build the data subsets in CAOS_{DS} is the random one, because it does not require the original classes of the instances, it is not significantly different in terms of accuracy than the other two strategies and it has a high speedup.

5.4. Discussion of CAOS and MOCK Results

The proposed approach has demonstrated that it is able to scale-up CAOS without considerably penalizing the accuracy and it has been compared with respect to MOCK (Handl and Knowles, 2007) in order to emphasize the results. MOCK is also based on the PESA-II algorithm but it uses the locus-based adjacency representation proposed in (Park and Song, 1998) for representing the individuals. That is, each individual encodes a reflexive directed unlabeled graph that connects pairs of examples using an integer encoding scheme. More specifically, each individual consists of m genes $\{x_1, x_2, \dots, x_m\}$, where m is the number of examples of the training data set and x_i ranges in $[1, m]$. Thence, each gene x_i indicates that there exists an arrow connecting instance i with instance x_i . It can be observed that this representation cannot properly scale-up the memory usage when the algorithm is applied to large data because the individual size depends on the number of instances. Thus, this representation does not allow the use of data subsets due to the fact that all the instances are needed for building the individuals.

The accuracy of the best results of MOCK, CAOS_{CD} and the three CAOS_{DS} strategies of the real-world data sets is analyzed using the Holm's test and illustrated in Table 4(a) For carrying out this analysis, four data sets (*letter-recognition*, *magic*, *pendigits* and *shuttle*) where MOCK was not able to obtain the results with a reasonable memory usage are not considered. Analyzing the results it can be observed that MOCK is significantly worse than CAOS_{CD} and than all the CAOS_{DS} strategies when a 20% of information or more is used. When less of a 20% of instances are used, MOCK is worse but not significantly different than CAOS_{DS} strategies. Finally, MOCK has a lower computational time than CAOS_{CD} for the different data sets as Table 11 and Table 12 show in the supplementary material. However, CAOS_{DS} strategies are significantly better than MOCK in terms of computational time independently of the size of the strata as Table 4(b) shows. Thus, the proposed data subset approach can obtain better results than MOCK in a lower computational time, and it can be affirmed that CAOS_{DS} is a promising scalable MC algorithm based on MOEAs.

Table 4. Comparison of (a) the accuracy and (b) the execution time of the algorithms in the real-world data sets (the four data sets where MOCK cannot obtain results are excluded from the analysis, so 26 data sets are used) using the post-hoc Holm's procedure with $\alpha = 0.05$. The algorithms compared are MOCK and CAOS using the CAOS_{CD} and the three CAOS_{DS} strategies to generate data subsets: based on classes (CAOS_{DS} Classes), random (CAOS_{DS} Random) and based on clusters (CAOS_{DS} Clusters). The results are showed for 34%, 10% and 4% of information used from the complete data set. The symbols \oplus and \ominus show that the method in the row obtained results that were significantly higher/lower than those obtained with the method in the column. Similarly, the symbols + and - denote a non-significant higher/lower results. Symbol = indicates that the algorithms have the same rank. The Iman and Davenport statistic is calculated according to F-distribution with 4 and 100 degrees of freedom. The results with the rest of % of instances are detailed in Table 9 and 10 in supplementary material.

(a)								
Instances	Strategies	MOCK	CAOS _{CD}	CAOS _{DS} Classes	CAOS _{DS} Random	CAOS _{DS} Clusters	Average Rank	p-value
34%	MOCK		\ominus	\ominus	\ominus	\ominus	4.346	$3.449E - 7$
	CAOS _{CD}	\oplus		\oplus	\oplus	\oplus	1.962	
	CAOS _{DS} Classes	\oplus	-		-	-	3.077	
	CAOS _{DS} Random	\oplus	-	\oplus		\oplus	2.769	
	CAOS _{DS} Clusters	\oplus	-	\oplus	-		2.846	
10%	MOCK		\ominus	-	-	\ominus	4.269	$1.108E - 9$
	CAOS _{CD}	\oplus		\oplus	\oplus	\oplus	1.577	
	CAOS _{DS} Classes	\oplus	\ominus		=	-	3.192	
	CAOS _{DS} Random	\oplus	\ominus	=		-	3.192	
	CAOS _{DS} Clusters	\oplus	\ominus	\oplus	\oplus		2.769	
4%	MOCK		\ominus	-	-	-	4.000	$4.580E - 8$
	CAOS _{CD}	\oplus		\oplus	\oplus	\oplus	1.538	
	CAOS _{DS} Classes	\oplus	\ominus		\oplus	\oplus	3.038	
	CAOS _{DS} Random	\oplus	\ominus	-		-	3.346	
	CAOS _{DS} Clusters	\oplus	\ominus	-	\oplus		3.077	
(b)								
Instances	Strategies	MOCK	CAOS _{CD}	CAOS _{DS} Classes	CAOS _{DS} Random	CAOS _{DS} Clusters	Average Rank	p-value
34%	MOCK		\oplus	\ominus	\ominus	\ominus	1.731	0
	CAOS _{CD}	-		\ominus	\ominus	\ominus	1.308	
	CAOS _{DS} Classes	\oplus	\oplus		\oplus	\oplus	4.231	
	CAOS _{DS} Random	\oplus	\oplus	-		\oplus	3.961	
	CAOS _{DS} Clusters	\oplus	\oplus	-	-		3.769	
10%	MOCK		\oplus	\ominus	\ominus	\ominus	1.692	0
	CAOS _{CD}	-		\ominus	\ominus	\ominus	1.308	
	CAOS _{DS} Classes	\oplus	\oplus		\oplus	\oplus	4.154	
	CAOS _{DS} Random	\oplus	\oplus	-		\oplus	4.135	
	CAOS _{DS} Clusters	\oplus	\oplus	-	-		3.712	
4%	MOCK		\oplus	\ominus	\ominus	\ominus	1.692	0
	CAOS _{CD}	-		\ominus	\ominus	\ominus	1.308	
	CAOS _{DS} Classes	\oplus	\oplus		\oplus	\oplus	4.442	
	CAOS _{DS} Random	\oplus	\oplus	-		\oplus	4.192	
	CAOS _{DS} Clusters	\oplus	\oplus	-	-		3.365	

5.5. CAOS Results in Large and Medium Data Sets

The analysis of the performance of the three CAOS_{DS} strategies in artificial and real-world data sets concluded that the random strata strategy is the most suitable one. In order to analyze the potential of CAOS_{DS} with this strategy applied to data sets with a considerable number of instances, it was tested with 33 data sets considered of medium size and with 33 data sets considered of large size. The medium size data sets have a number of instances from 2000 to 8000, and the large size data sets have a number of instances from 13000 to 50000.

The experiments show that the speedup of both type of data sets is very high (see Figure 3(b)), and this represents an important improvement (e.g., in the shuttle data set, which has 50000 instances, CAOS_{CD} needs near 40 hours to obtain a clustering result, and CAOS_{DS} with a random strata strategy needs 10 hours using the 50% of instances and 3 hours using the 25% of them, as Table 11 in supplementary material shows). It is obvious that if less data is considered, the accuracy results will be worse, as statistical tests show in Subsection 5.3. However, the loss of accuracy is not very considerable when more than the 20% of instances are used (see Figure 3(a)). Thus, the results

obtained when the 25% of the instances are used in each cycle of the EA are analyzed to better illustrate this issue because it is considered (according to Figure 3) that it is a configuration with a good trade-off between accuracy and execution time.

Figure 4 and Figure 5 show the accuracy result and the overall execution time of CAOS_{CD} and the random strata strategy of CAOS_{DS} for each one of the large and medium-sized data sets. Both figures show the average results of the 20 runs done. It can be observed that the accuracy results are preserved or slightly penalized but, on the other hand, the execution time is strongly improved. Moreover, it is important to emphasize that the main idea of CAOS_{DS} is to considerably improve the performance of the algorithm with small impact on the accuracy, and this can be clearly observed in data sets of medium and large size.

6. Conclusions and Further Work

MC based on EA is a data mining technique focused on identifying data relationships according to multiple criteria to properly understand huge and complex databases. Although its search capability outlines from the rest of similar techniques, its main lack is the high cost in terms of computational time and memory usage when it is applied to a large data sets.

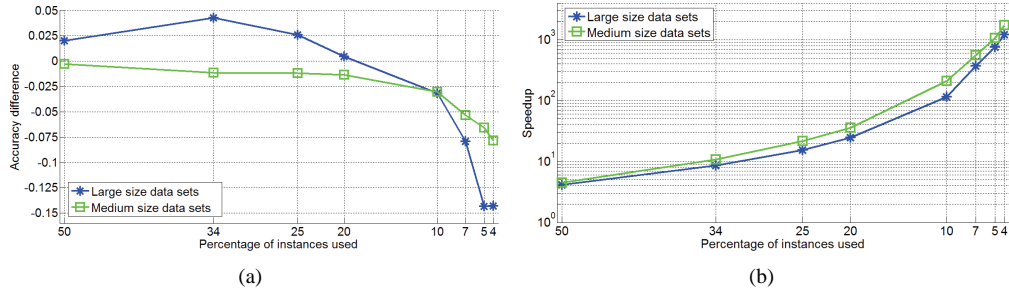


Fig. 3. Accuracy difference (a) and overall speedup (b) of the CAOS_{DS} random strata strategy in large and medium size data sets.

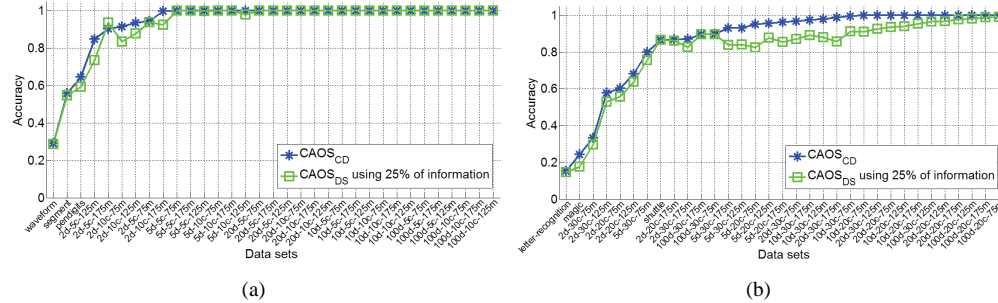


Fig. 4. Average accuracy of 20 runs of CAOS_{CD} and CAOS_{DS} using 25% of information in each generation for (a) medium and (b) large data sets.

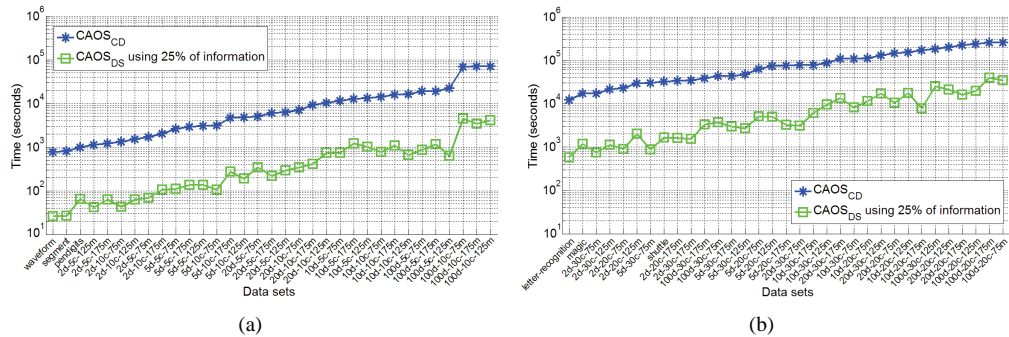


Fig. 5. Average of the overall clustering time of 20 runs of CAOS_{CD} and CAOS_{DS} using the 25% of information in each generation for (a) medium and (b) large size data sets.

This work has presented an approach focused on reducing the impact of the volume of data in the EA by means of the stratification of the complete data set into disjoint strata and alternate them following a Round Robin policy in each cycle of the genetic algorithm. More specifically, a supervised and two unsupervised stratification strategies are proposed and their performance is analyzed using 106 real and artificial data sets in the CAOS algorithm: (1) according to the original classes of the data set, (2) selecting random instances from the data set, and (3) according to the clusters found applying a fast method called subtractive clustering.

The experimentation showed that the speedup of the three strategies is very high, and this considerably improves the computational performance of the system. Moreover, it can be observed that the two unsupervised strategies used to build the strata are not significantly different from the strategy based on classes in terms of accuracy so they can be considered equivalent to the strategy based on classes in terms of accuracy. Furthermore, the strategy based on random strata has a higher speedup than the cluster based strategy because the last one needs to build approximative clusters at the begin of the algorithm and this has a high cost with very large data sets. Thus, the random stratum strategy is the most suitable to scaling-up CAOS. On the other hand, there are statistically significant differences among these three strategies and the approach based on using the complete data set even the accuracy differences among them are relatively small while they considerably reduce the computational time of the algorithm. It is important to highlight that the size of each stratum affects the performance of the CAOS algorithm. The computational time of CAOS decreases as much smaller is the size of the strata but it is more difficult to obtain consistent stratum according to the original data set so it affects the accuracy of the system. The last part of the experimentation has also compared the results of the proposed strategies with respect to MOCK and they have obtained results significantly better in terms of computational time and accuracy. Thus, it can be concluded that the proposed approach is a promising scalable MC algorithm based on MOEAs and it can be properly applied to large data. Moreover, this approach can also be applied in situations where losing some accuracy can be accepted if it is possible to obtain results in a reasonable time.

Finally, this analysis has set the basis for further conducting research on multiobjective evolutionary clustering applied to large data sets in two research topics. First, the consequences of applying stratification methods for scaling-up CAOS with other individual representations. Second, the application of other data mining techniques for large data sets (Bacardit and Llorà, 2009) such as Parallel EAs.

References

- Bacardit, J., 2004. Pittsburgh Genetic-Based Machine Learning in the Data Mining era: Representations, generalization, and run-time. Ph.D. thesis. Enginyeria i Arquitectura La Salle, Universitat Ramon Llull.
- Bacardit, J., Llorà, X., 2009. Large scale data mining using genetics-based machine learning, in: Proceedings of the 11th Annual Conference Companion on GECCO: Late Breaking Papers, ACM, pp. 3381–3412.
- Cano, J.R., García, S., Herrera, F., 2008. Subgroup discovery in large size data sets preprocessed using stratified instance selection for increasing the presence of minority classes. *Pattern Recognition Letters* 29, 2156–2164.
- Cano, J.R., Herrera, F., Lozano, M., 2006. On the combination of evolutionary algorithms and stratified strategies for training set selection in data mining. *Applied Soft Computing* 2006, 323–332.
- Cantu-Paz, E., 2000. *Efficient and Accurate Parallel Genetic Algorithms*. Kluwer Academic Publishers, Norwell, MA, USA.
- Chiu, S.L., 1994. Fuzzy model identification based on cluster estimation. *Journal of Intelligent and Fuzzy Systems*.
- Coello, C.A., 1999. A comprehensive survey of evolutionary multiobjective optimization techniques. *Knowledge and Information Systems* 1, 269–308.
- Corne, D.W., Jerram, N.R., Knowles, J.D., Oates, M.J., 2001. PESA-II: Region-based selection in evolutionary multiobjective optimization, in: Proceedings of the Genetic and Evolutionary Computation Conference, Morgan Kaufmann Publishers. pp. 283–290.
- Demšar, J., 2006. Statistical comparisons of classifiers over multiple data sets. *Journal of Machine Learning Research* 7, 1–30.
- Derrac, J., García, S., Herrera, F., 2010. Stratified prototype selection based on a steady-state memetic algorithm: A study of scalability. *Memetic Computing* 2, 183–199.
- Frank, A., Asuncion, A., 2010. UCI machine learning repository. URL: <http://archive.ics.uci.edu/ml>.
- Freitas, A.A., 2002. *Data Mining and Knowledge Discovery with Evolutionary Algorithms*. Springer-Verlag New York, Inc., Secaucus, NJ, USA.
- Friedman, M., 1940. A comparison of alternative tests of significance for the problem of m rankings. *Annals of Mathematical Statistics* 11, 86–92.
- García, S., Herrera, F., 2008. An extension on statistical comparisons of classifiers over multiple data sets for all pairwise comparisons. *Journal of Machine Learning Research* 9, 2677–2694.
- García-Piquer, A., Fornells, A., Bacardit, J., Orriols-Puig, A., Golobardes, E., 2013. Large-scale experimental evaluation of cluster representations for multiobjective evolutionary clustering. *IEEE Transactions on Evolutionary Computation*, 36–53.
- García-Piquer, A., Fornells, A., Orriols-Puig, A., Corral, G., Golobardes, E., 2012. Data Classification through an Evolutionary Approach Based on Multiple Criteria. *Knowledge and Information Systems* 33, 35–56.
- Goldberg, D.E., 2002. *The Design of Innovation*. Kluwer Academic Publishers, Massachusetts, US.
- Halkidi, M., Batistakis, Y., Vazirgiannis, M., 2001. On clustering validation techniques. *Journal of Intelligent Information Systems* 17, 107–145.
- Handl, J., Knowles, J., 2007. An evolutionary approach to multiobjective clustering. *IEEE Transactions on Evolutionary Computation* 1, 56–76.
- Holm, S., 1979. A simple sequentially rejective multiple test procedure. *Scandinavian Journal of Statistics* 6, 65–70.
- Hore, P., Hall, L.O., Goldgof, D.B., 2009. A scalable framework for cluster ensembles. *Pattern Recognition* 42, 676–688.
- Hruschka, E.R., Campello, R.J.G.B., Freitas, A.A., de Carvalho, A.C.P.L.F., 2009. A survey of evolutionary algorithms for clustering. *IEEE Transactions on Systems, Man and Cybernetics, Applications and Reviews* 39, 133–155.
- Iredi, S., Merkle, D., Middendorf, M., 2000. Bi-criterion optimization with multi colony ant algorithms, in: Proceedings of the 1st International Conference on Evolutionary Multi-Criterion Optimization, Springer. pp. 359–372.
- Kargupta, H., Han, J., Yu, P.S., Motwani, R., Kumar, V., 2009. *Next Generation of Data Mining*. Chapman & Hall/CRC data mining and knowledge discovery series, CRC Press, USA.
- Kaufman, L., Rousseeuw, P.J., 1990. *Finding groups in data: An introduction to cluster analysis*. John Wiley & Sons.
- Park, Y., Song, M., 1998. A genetic algorithm for clustering problems, in: Proceedings of the 3rd Annual Conference on Genetic Programming, Morgan Kaufmann. pp. 568–575.
- Saha, S., Bandyopadhyay, S., 2010. A new multiobjective clustering technique based on the concepts of stability and symmetry. *Knowledge and Information Systems* 23, 1–27.
- Yeung, K., Ruzzo, W., 2001. Details of the adjusted rand index and clustering algorithms. supplement to the paper “an empirical study on principal component analysis for clustering gene expression data”. *Science* 17, 763–774.
- Zitzler, E., Deb, K., Thiele, L., 2000. Comparison of multiobjective evolutionary algorithms. *Evolutionary computation* 8, 173–195.

Supplementary Material

Table 5. Summary of the characteristics of the 30 real-world data sets used, sorted by their number of instances. The columns are referred to the number of instances (nI), to the number of attributes (nA) and to the number of classes (nC).

Data set	nI	nA	nC	Data set	nI	nA	nC
shuttle	58000	9	7	crx	690	15	2
letter-recognition	20000	16	26	balance	625	4	3
magic	19020	10	2	wdbc	569	30	2
pendigits	7494	17	10	housevotes	435	16	2
waveform	5000	40	3	dermatology	366	35	6
segment	2310	19	7	ionosphere	351	34	2
yeast	1484	9	10	bpa	345	6	2
contraceptives	1473	9	3	liver-disorders	345	6	2
biopsia	1027	24	2	ecoli	336	8	8
vowel	990	13	11	heart-statlog	270	13	2
mammographic	961	5	2	thyroids	215	5	2
vehicle	846	18	4	glass	214	9	6
pim	768	8	2	sonar	208	60	2
transfusion	748	4	2	wpbc	198	33	2
wisconsin	699	9	2	iris	150	4	3

Table 6. Summary of the characteristics of the 75 artificial data sets used, sorted by their number of instances. The columns are referred to the number of instances (nI), to the number of attributes (nA) and to the number of classes (nC).

Data set	nI	nA	nC	Data set	nI	nA	nC
10d-30c-175m	23898	10	30	100d-10c-125m	7021	100	10
10d-30c-75m	23471	10	30	5d-10c-175m	6981	5	10
5d-30c-75m	23234	5	30	20d-10c-75m	6848	20	10
100d-30c-175m	22788	100	30	2d-10c-75m	6798	2	10
20d-30c-75m	22470	20	30	2d-10c-175m	6586	2	10
2d-30c-175m	22229	2	30	10d-10c-75m	6513	10	10
5d-30c-125m	22038	5	30	20d-10c-125m	6084	20	10
10d-30c-125m	21974	10	30	20d-5c-125m	4473	20	5
2d-30c-125m	21846	2	30	2d-5c-75m	4306	2	5
20d-30c-175m	21491	20	30	20d-5c-75m	3877	20	5
5d-30c-175m	21129	5	30	100d-5c-175m	3870	100	5
20d-30c-125m	20986	20	30	10d-5c-75m	3781	10	5
100d-30c-75m	20561	100	30	10d-5c-175m	3664	10	5
2d-30c-75m	20370	2	30	5d-5c-175m	3575	5	5
100d-30c-125m	20156	100	30	2d-5c-175m	3564	2	5
2d-20c-125m	16097	2	20	2d-5c-125m	3525	2	5
5d-20c-175m	15675	5	20	100d-5c-75m	3445	100	5
20d-20c-125m	15508	20	20	10d-5c-125m	3382	10	5
2d-20c-75m	15012	2	20	20d-5c-175m	3272	20	5
20d-20c-175m	14970	20	20	5d-5c-75m	3184	5	5
10d-20c-75m	14830	10	20	5d-5c-125m	3124	5	5
20d-20c-75m	14491	20	20	100d-5c-125m	2955	100	5
5d-20c-125m	14261	5	20	20d-2c-75m	1796	20	2
10d-20c-175m	14023	10	20	20d-2c-125m	1789	20	2
10d-20c-125m	13875	10	20	100d-2c-75m	1763	100	2
100d-20c-75m	13790	100	20	20d-2c-175m	1685	20	2
100d-20c-125m	13702	100	20	5d-2c-125m	1636	5	2
100d-20c-175m	13421	100	20	2d-2c-75m	1560	2	2
2d-20c-175m	13355	2	20	2d-2c-125m	1497	2	2
5d-20c-75m	13289	5	20	10d-2c-125m	1418	10	2
5d-10c-75m	7875	5	10	5d-2c-75m	1352	5	2
10d-10c-125m	7609	10	10	2d-2c-175m	1245	2	2
100d-10c-75m	7394	100	10	100d-2c-125m	1146	100	2
20d-10c-175m	7200	20	10	100d-2c-175m	1132	100	2
2d-10c-125m	7078	2	10	10d-2c-175m	1132	10	2
100d-10c-175m	7071	100	10	5d-2c-175m	1042	5	2
10d-10c-175m	7070	10	10	10d-2c-75m	884	10	2
5d-10c-125m	7052	5	10				

Table 7. Comparison of the accuracy of the algorithms in the overall data sets (artificial and real-world) using the post-hoc Holm's procedure with $\alpha = 0.05$. The algorithms compared are CAOS using the CAOS_{CD} and the three CAOS_{DS} strategies to generate data subsets: based on classes (CAOS_{DS} Classes), random (CAOS_{DS} Random) and based on clusters (CAOS_{DS} Clusters). The results are showed for 50%, 34%, 25%, 20%, 10%, 7% and 4% of information used from the complete data set. The symbols \oplus and \ominus show that the method in the row obtained results that were significantly higher/lower than those obtained with the method in the column. Similarly, the symbols + and - denote a non-significant higher/lower results. The Iman and Davenport statistic is calculated according to F-distribution with 3 and 312 degrees of freedom.

Instances	Strategies	CAOS _{CD}	CAOS _{DS} Classes	CAOS _{DS} Random	CAOS _{DS} Cluster	Average Rank	p-value
50%	CAOS _{CD}		\oplus	\oplus	\oplus	1.957	$3.133E - 6$
	CAOS _{DS} Classes	\ominus		-	+	2.601	
	CAOS _{DS} Random	\ominus	+		+	2.591	
	CAOS _{DS} Clusters	\ominus	-	-		2.851	
34%	CAOS _{CD}		\oplus	\oplus	\oplus	1.856	$5.972E - 8$
	CAOS _{DS} Classes	\ominus		-	-	2.736	
	CAOS _{DS} Random	\ominus	+		+	2.678	
	CAOS _{DS} Clusters	\ominus	+	-		2.731	
25%	CAOS _{CD}		\oplus	\oplus	\oplus	1.827	$3.360E - 9$
	CAOS _{DS} Classes	\ominus		-	-	2.875	
	CAOS _{DS} Random	\ominus	+		+	2.625	
	CAOS _{DS} Clusters	\ominus	+	-		2.673	
20%	CAOS _{CD}		\oplus	\oplus	\oplus	1.745	$4.622E - 11$
	CAOS _{DS} Classes	\ominus		+	+	2.726	
	CAOS _{DS} Random	\ominus	-		+	2.760	
	CAOS _{DS} Clusters	\ominus	-	-		2.770	
10%	CAOS _{CD}		\oplus	\oplus	\oplus	1.673	$1.105E - 13$
	CAOS _{DS} Classes	\ominus		+	-	2.788	
	CAOS _{DS} Random	\ominus	-		-	2.851	
	CAOS _{DS} Clusters	\ominus	+	+		2.687	
7%	CAOS _{CD}		\oplus	\oplus	\oplus	1.644	$1.010E - 14$
	CAOS _{DS} Classes	\ominus		+	-	2.745	
	CAOS _{DS} Random	\ominus	-		-	2.875	
	CAOS _{DS} Clusters	\ominus	+	+		2.730	
4%	CAOS _{CD}		\oplus	\oplus	\oplus	1.596	0
	CAOS _{DS} Classes	\ominus		+	+	2.663	
	CAOS _{DS} Random	\ominus	-		-	2.961	
	CAOS _{DS} Clusters	\ominus	-	+		2.779	

Table 8. Comparison of the accuracy of the algorithms in the (a) artificial and (b) real data sets using the post-hoc Holm’s procedure with $\alpha = 0.05$. The algorithms compared are CAOS using the CAOS_{CD} and the three CAOS_{DS} strategies to generate data subsets: based on classes (CAOS_{DS} Classes), random (CAOS_{DS} Random) and based on clusters (CAOS_{DS} Clusters). The results are showed for 50%, 34%, 25%, 20%, 10%, 7% and 4% of information used from the complete data set. The symbols \oplus and \ominus show that the method in the row obtained results that were significantly higher/lower than those obtained with the method in the column. Similarly, the symbols + and – denote a non-significant higher/lower results. The Iman and Davenport statistic is calculated according to F-distribution with 3 and 222 degrees of freedom.

(a)							
Instances	Strategies	CAOS _{CD}	CAOS _{DS} Classes	CAOS _{DS} Random	CAOS _{DS} Clusters	Average rank	p-value
50%	CAOS _{CD}		\oplus	\oplus	\oplus	1.946	$5.131E - 5$
	CAOS _{DS} Classes	\ominus		–	+	2.621	
	CAOS _{DS} Random	\ominus	+		+	2.527	
	CAOS _{DS} Clusters	\ominus	–	–		2.905	
34%	CAOS _{CD}		\oplus	\oplus	\oplus	1.871	$1.705E - 5$
	CAOS _{DS} Classes	\ominus		+	+	2.662	
	CAOS _{DS} Random	\ominus	–		+	2.689	
	CAOS _{DS} Clusters	\ominus	–	–		2.777	
25%	CAOS _{CD}		\oplus	\oplus	\oplus	1.818	$1.818E - 6$
	CAOS _{DS} Classes	\ominus		+	+	2.777	
	CAOS _{DS} Random	\ominus	–		–	2.669	
	CAOS _{DS} Clusters	\ominus	–	+		2.736	
20%	CAOS _{CD}		\oplus	\oplus	\oplus	1.798	$5.494E - 7$
	CAOS _{DS} Classes	\ominus		+	–	2.770	
	CAOS _{DS} Random	\ominus	–		–	2.804	
	CAOS _{DS} Clusters	\ominus	+	+		2.628	
10%	CAOS _{CD}		\oplus	\oplus	\oplus	1.730	$2.634E - 8$
	CAOS _{DS} Classes	\ominus		+	+	2.716	
	CAOS _{DS} Random	\ominus	–		–	2.817	
	CAOS _{DS} Clusters	\ominus	–	+		2.736	
7%	CAOS _{CD}		\oplus	\oplus	\oplus	1.716	$7.542E - 9$
	CAOS _{DS} Classes	\ominus		+	–	2.757	
	CAOS _{DS} Random	\ominus	–		–	2.878	
	CAOS _{DS} Clusters	\ominus	+	+		2.649	
4%	CAOS _{CD}		\oplus	\oplus	\oplus	1.676	$4.689E - 10$
	CAOS _{DS} Classes	\ominus		+	+	2.662	
	CAOS _{DS} Random	\ominus	–		–	2.946	
	CAOS _{DS} Clusters	\ominus	–	+		2.716	

(b)							
Instances	Strategies	CAOS _{CD}	CAOS _{DS} Classes	CAOS _{DS} Random	CAOS _{DS} Clusters	Average rank	p-value
50%	CAOS _{CD}					1.983	2.377
	CAOS _{DS} Classes					2.550	
	CAOS _{DS} Random					2.750	
	CAOS _{DS} Clusters					2.716	
The Friedman’s test cannot reject the null hypothesis that all the learning algorithms obtain the same results on average.							
34%	CAOS _{CD}		\oplus	+	+	1.816	0.005
	CAOS _{DS} Classes	\ominus		–	–	2.917	
	CAOS _{DS} Random	–	+		+	2.650	
	CAOS _{DS} Clusters	–	–	+		2.617	
25%	CAOS _{CD}		\oplus	+	+	1.850	0.001
	CAOS _{DS} Classes	\ominus		–	–	3.117	
	CAOS _{DS} Random	–	+		+	2.517	
	CAOS _{DS} Clusters	–	–	+		2.519	
20%	CAOS _{CD}		\oplus	\oplus	\oplus	1.617	$2.581E - 5$
	CAOS _{DS} Classes	\ominus		+	+	2.617	
	CAOS _{DS} Random	\ominus	–		+	2.650	
	CAOS _{DS} Clusters	\ominus	–	–		3.117	
10%	CAOS _{CD}		\oplus	\oplus	\oplus	1.533	$4.853E - 6$
	CAOS _{DS} Classes	\ominus		–	–	2.966	
	CAOS _{DS} Random	\ominus	+		–	2.933	
	CAOS _{DS} Clusters	\ominus	+	+		2.567	
7%	CAOS _{CD}		\oplus	\oplus	\oplus	1.467	$1.502E - 6$
	CAOS _{DS} Classes	\ominus		+	+	2.733	
	CAOS _{DS} Random	\ominus	–		+	2.867	
	CAOS _{DS} Clusters	\ominus	–	–		2.933	
4%	CAOS _{CD}		\oplus	\oplus	\oplus	1.400	$8.677E - 8$
	CAOS _{DS} Classes	\ominus		+	+	2.667	
	CAOS _{DS} Random	\ominus	–		–	3.000	
	CAOS _{DS} Clusters	\ominus	–	+		2.933	

Table 9. Comparison of the accuracy of the algorithms in the real-world data sets (the four data sets where MOCK cannot obtain results are excluded from the analysis, so 26 data sets are used) using the post-hoc Holm’s procedure with $\alpha = 0.05$. The algorithms compared are MOCK and CAOS using the CAOS_{CD} and the three CAOS_{DS} strategies to generate data subsets: based on classes (CAOS_{DS} Classes), random (CAOS_{DS} Random) and based on clusters (CAOS_{DS} Clusters). The results are showed for 50%, 34%, 25%, 20%, 10%, 7% and 4% of information used from the complete data set. The symbols \oplus and \ominus show that the method in the row obtained results that were significantly higher/lower than those obtained with the method in the column. Similarly, the symbols + and – denote a non-significant higher/lower results. Symbol = indicates that the algorithms have the same rank. The Iman and Davenport statistic is calculated according to F-distribution with 4 and 100 degrees of freedom.

Instances	Strategies	MOCK	CAOS _{CD}	CAOS _{DS} Classes	CAOS _{DS} Random	CAOS _{DS} Clusters	Average rank	p-value
50%	MOCK		\ominus	\ominus	\ominus	\ominus	4.462	$2.996E - 7$
	CAOS _{CD}	\oplus		+	+	+	2.154	
	CAOS _{DS} Classes	\oplus	–		+	+	2.692	
	CAOS _{DS} Random	\oplus	–	–		=	2.846	
	CAOS _{DS} Clusters	\oplus	–	–	=		2.846	
34%	MOCK		\ominus	\ominus	\ominus	\ominus	4.346	$3.449E - 7$
	CAOS _{CD}	\oplus		+	+	+	1.962	
	CAOS _{DS} Classes	\oplus	–		–	–	3.077	
	CAOS _{DS} Random	\oplus	–	+		+	2.769	
	CAOS _{DS} Clusters	\oplus	–	+	–		2.846	
25%	MOCK		\ominus	\ominus	\ominus	\ominus	4.385	$4.880E - 8$
	CAOS _{CD}	\oplus		\oplus	+	+	2.000	
	CAOS _{DS} Classes	\oplus	\ominus		–	–	3.308	
	CAOS _{DS} Random	\oplus	–	+		–	2.692	
	CAOS _{DS} Clusters	\oplus	–	+	+		2.615	
20%	MOCK		\ominus	\ominus	\ominus	\ominus	4.423	$4.378E - 10$
	CAOS _{CD}	\oplus		+	+	+	1.731	
	CAOS _{DS} Classes	\oplus	–		–	+	2.769	
	CAOS _{DS} Random	\oplus	–	+		+	2.731	
	CAOS _{DS} Clusters	\oplus	\ominus	–	–		3.346	
10%	MOCK		\ominus	–	–	\ominus	4.269	$1.108E - 9$
	CAOS _{CD}	\oplus		\oplus	\oplus	\oplus	1.577	
	CAOS _{DS} Classes	+	\ominus		=	–	3.192	
	CAOS _{DS} Random	+	\ominus	=		–	3.192	
	CAOS _{DS} Clusters	\oplus	\ominus	+	+		2.769	
7%	MOCK		\ominus	–	–	–	4.115	$4.816E - 9$
	CAOS _{CD}	\oplus		\oplus	\oplus	\oplus	1.500	
	CAOS _{DS} Classes	+	\ominus		+	+	3.000	
	CAOS _{DS} Random	+	\ominus	–		+	3.154	
	CAOS _{DS} Clusters	+	\ominus	–	–		3.231	
4%	MOCK		\ominus	–	–	–	4.000	$4.580E - 8$
	CAOS _{CD}	\oplus		\oplus	\oplus	\oplus	1.538	
	CAOS _{DS} Classes	+	\ominus		+	+	3.038	
	CAOS _{DS} Random	+	\ominus	–		–	3.346	
	CAOS _{DS} Clusters	+	\ominus	–	+		3.077	

Table 10. Comparison of the execution time of the algorithms in the real-world data sets (the four data sets where MOCK cannot obtain results are excluded from the analysis, so 26 data sets are used) using the post-hoc Holm’s procedure with $\alpha = 0.05$. The algorithms compared are MOCK and CAOS using the CAOS_{CD} and the three CAOS_{DS} strategies to generate data subsets: based on classes (CAOS_{DS} Classes), random (CAOS_{DS} Random) and based on clusters (CAOS_{DS} Clusters). The results are showed for 50%, 34%, 25%, 20%, 10%, 7% and 4% of information used from the complete data set. The symbols \oplus and \ominus show that the method in the row obtained results that were significantly higher/lower than those obtained with the method in the column. Similarly, the symbols + and – denote a non-significant higher/lower results. Symbol = indicates that the algorithms have the same rank. It must be worth noting that, in terms of execution time, high ranks are better due to low values of time are required. The Iman and Davenport statistic is calculated according to F-distribution with 4 and 100 degrees of freedom.

Instances	Strategies	MOCK	CAOS _{CD}	CAOS _{DS} Classes	CAOS _{DS} Random	CAOS _{DS} Clusters	Average rank	p-value
50%	MOCK		+	\ominus	\ominus	\ominus	1.923	0
	CAOS _{CD}	–		\ominus	\ominus	\ominus	1.308	
	CAOS _{DS} Classes	\oplus	\oplus	=	=	=	3.923	
	CAOS _{DS} Random	\oplus	\oplus	=	=	=	3.923	
	CAOS _{DS} Clusters	\oplus	\oplus	=	=	=	3.923	
34%	MOCK		+	\ominus	\ominus	\ominus	1.731	0
	CAOS _{CD}	–		\ominus	\ominus	\ominus	1.308	
	CAOS _{DS} Classes	\oplus	\oplus		+	+	4.231	
	CAOS _{DS} Random	\oplus	\oplus	–		+	3.961	
	CAOS _{DS} Clusters	\oplus	\oplus	–	–		3.769	
25%	MOCK		+	\ominus	\ominus	\ominus	1.692	0
	CAOS _{CD}	–		\ominus	\ominus	\ominus	1.308	
	CAOS _{DS} Classes	\oplus	\oplus		+	+	4.250	
	CAOS _{DS} Random	\oplus	\oplus	–		–	3.846	
	CAOS _{DS} Clusters	\oplus	\oplus	–	+		3.904	
20%	MOCK		+	\ominus	\ominus	\ominus	1.692	0
	CAOS _{CD}	–		\ominus	\ominus	\ominus	1.308	
	CAOS _{DS} Classes	\oplus	\oplus		+	+	4.154	
	CAOS _{DS} Random	\oplus	\oplus	–		+	4.039	
	CAOS _{DS} Clusters	\oplus	\oplus	–	–		3.808	
10%	MOCK		+	\ominus	\ominus	\ominus	1.692	0
	CAOS _{CD}	–		\ominus	\ominus	\ominus	1.308	
	CAOS _{DS} Classes	\oplus	\oplus		+	+	4.154	
	CAOS _{DS} Random	\oplus	\oplus	–		+	4.135	
	CAOS _{DS} Clusters	\oplus	\oplus	–	–		3.712	
7%	MOCK		+	\ominus	\ominus	\ominus	1.692	0
	CAOS _{CD}	–		\ominus	\ominus	\ominus	1.308	
	CAOS _{DS} Classes	\oplus	\oplus		+	+	4.154	
	CAOS _{DS} Random	\oplus	\oplus	–		+	4.135	
	CAOS _{DS} Clusters	\oplus	\oplus	–	–		3.712	
4%	MOCK		+	\ominus	\ominus	\ominus	1.692	0
	CAOS _{CD}	–		\ominus	\ominus	\ominus	1.308	
	CAOS _{DS} Classes	\oplus	\oplus		+	+	4.442	
	CAOS _{DS} Random	\oplus	\oplus	–		+	4.192	
	CAOS _{DS} Clusters	\oplus	\oplus	–	–		3.365	

Table 11. Execution time in seconds of CAOS using the CAOS_{CD} and the three CAOS_{DS} strategies to generate data subsets: based on classes (CAOS_{DS} Classes), random (CAOS_{DS} Random) and based on clusters (CAOS_{DS} Clusters). The results are showed for 50%, 34%, 25% and 20% of information used from the complete data set. Moreover, the execution time of MOCK algorithm is presented. The – symbol indicates that the corresponding algorithm was not able to obtain the results with a reasonable memory usage.

Data set	MOCK	CAOS _{CD}	CAOS _{DS} (50%)			CAOS _{DS} (34%)			CAOS _{DS} (25%)			CAOS _{DS} (20%)		
			Classes	Random	Clusters	Classes	Random	Clusters	Classes	Random	Clusters	Classes	Random	Clusters
balance	43.70	78.74	16.74	17.72	13.43	3.95	4.02	5.02	2.48	2.54	2.50	1.28	1.10	1.27
biopsia	176.65	511.88	108.58	135.70	168.29	48.01	62.10	44.81	22.64	26.98	22.30	10.99	15.72	11.62
bpa	21.52	27.37	5.42	5.94	5.04	1.72	2.58	2.10	0.92	0.92	1.18	0.50	0.68	0.66
contraceptives	142.04	349.84	72.16	90.85	85.72	34.39	32.51	28.87	14.24	16.39	17.74	9.65	11.11	9.13
crx	144.80	23.04	6.67	5.35	6.51	3.78	3.52	2.65	1.55	1.83	1.79	1.09	1.12	1.06
dermatology	52.52	146.08	29.58	30.10	32.80	9.82	10.81	10.11	4.78	7.30	4.43	2.11	1.66	1.95
ecoli	23.13	29.38	5.61	5.32	5.23	1.54	1.72	1.71	1.03	1.01	0.70	0.45	0.52	0.75
glass	16.09	12.77	2.27	2.74	2.72	0.92	1.10	1.20	0.59	0.52	0.39	0.50	0.32	0.49
heart-statlog	28.94	30.02	7.11	4.95	6.34	2.01	2.70	1.82	1.27	1.72	1.07	0.68	0.53	0.72
housevotes	42.72	9.75	0.34	0.35	3.16	1.30	1.29	1.52	0.71	0.65	0.84	0.53	0.73	0.52
ionosphere	57.85	118.16	22.22	26.06	29.79	7.17	9.08	6.23	3.88	3.44	4.87	2.21	3.03	2.76
iris	64.90	4.71	0.86	1.37	1.15	0.44	0.37	0.47	0.37	0.50	0.50	0.32	0.30	0.18
letter-recognition	–	46933.41	6611.63	8904.71	7571.54	5346.95	4250.12	5618.74	2844.15	2712.04	3716.18	2290.48	1626.48	2172.04
liver-disorders	22.83	35.19	5.82	6.97	4.16	2.34	2.62	1.82	0.96	1.04	1.05	0.49	0.62	0.69
magic	–	38392.73	15943.77	12473.66	14324.35	7643.63	6177.70	5608.94	2637.42	3344.59	2553.88	2184.74	1712.92	2735.09
mammography	60.78	10.58	2.22	2.33	2.55	1.19	1.26	1.17	0.80	0.56	0.82	0.47	0.48	0.52
pendigits	–	14089.59	4554.12	3573.42	4330.80	1584.21	1636.68	1484.44	914.70	763.11	873.11	403.51	419.59	483.48
pim	65.38	125.54	33.14	31.35	25.55	9.23	9.56	9.28	5.66	5.50	4.52	3.48	2.27	2.66
segment	279.67	986.58	317.80	282.72	332.02	135.50	124.62	144.94	51.26	64.40	65.80	43.36	37.77	47.40
shuttle	–	145769.48	39513.02	37278.12	40372.78	17929.51	19899.96	18259.26	10744.45	10444.68	10666.89	7297.98	7547.09	7615.68
sonar	52.86	125.57	25.91	22.54	24.61	5.91	9.63	7.98	4.25	4.78	4.05	3.26	2.46	3.10
thyroids	86.27	28.67	4.90	5.24	4.83	2.59	2.63	2.89	11.59	18.99	10.54	6.94	4.03	8.41
transfusion	36.14	19.28	4.85	4.45	6.05	2.72	3.03	3.08	1.66	1.63	1.68	1.05	0.78	1.24
vehicle	97.22	358.72	67.10	65.15	57.30	23.36	19.51	21.03	9.21	12.57	9.89	6.58	5.49	7.46
vowel	100.60	231.39	47.88	56.15	59.28	18.33	17.86	20.44	8.58	9.76	10.91	4.65	6.41	5.16
waveform	1437.58	22075.68	4287.51	4289.13	4198.36	1408.31	1258.11	1513.80	666.65	630.38	704.95	344.98	392.73	375.29
wdbc	90.45	264.08	43.03	40.76	52.11	17.37	14.86	17.48	6.92	9.28	7.29	3.82	5.21	6.85
wisc	45.49	34.54	13.50	9.83	7.60	4.77	4.68	6.31	2.19	2.15	2.30	1.94	2.05	1.10
wpbc	26.72	63.19	15.24	11.36	14.05	3.71	4.78	5.38	3.79	2.73	2.24	1.54	1.83	1.75
yeast	153.37	400.10	83.50	74.58	82.52	30.89	26.39	26.46	11.61	13.74	12.61	8.62	8.32	8.86

Table 12. Execution time in seconds of CAOS using the CAOS_{CD} and the three CAOS_{DS} strategies to generate data subsets: based on classes (CAOS_{DS} Classes), random (CAOS_{DS} Random) and based on clusters (CAOS_{DS} Clusters). The results are showed for 10%, 7%, 5% and 4% of information used from the complete data set. Moreover, the execution time of MOCK algorithm is presented. The – symbol indicates that the corresponding algorithm was not able to obtain the results with a reasonable memory usage.

Data set	MOCK	CAOS _{CD}	CAOS _{DS} (10%)			CAOS _{DS} (7%)			CAOS _{DS} (5%)			CAOS _{DS} (4%)		
			Classes	Random	Clusters	Classes	Random	Clusters	Classes	Random	Clusters	Classes	Random	Clusters
balance	43.70	78.74	0.36	0.37	0.38	0.15	0.13	0.20	0.12	0.06	0.18	0.05	0.06	0.15
biopsia	176.65	511.88	2.78	2.21	1.95	0.75	0.84	1.56	0.63	0.46	1.05	0.48	0.58	0.70
bpa	21.52	27.37	0.16	0.12	0.15	0.12	0.12	0.11	0.08	0.12	0.07	0.08	0.07	0.09
contraceptives	142.04	349.84	1.72	1.89	2.62	0.61	1.12	0.83	0.43	0.42	0.78	0.26	0.35	0.66
crx	144.80	23.04	0.33	0.32	0.42	0.19	0.19	0.29	0.15	0.15	0.29	0.15	0.14	0.27
dermatology	52.52	146.08	0.62	0.68	0.81	0.37	0.39	0.50	0.25	0.29	0.41	0.27	0.29	0.27
ecoli	23.13	29.38	0.17	0.16	0.16	0.13	0.12	0.11	0.12	0.10	0.10	0.07	0.09	0.10
glass	16.09	12.77	0.13	0.12	0.12	0.13	0.11	0.10	0.08	0.10	0.10	0.07	0.12	0.07
heart-statlog	28.94	30.02	0.23	0.24	0.28	0.13	0.14	0.15	0.10	0.12	0.11	0.11	0.11	0.12
housevotes	42.72	9.75	0.20	0.20	0.28	0.18	0.14	0.22	0.14	0.16	0.23	0.14	0.14	0.20
ionosphere	57.85	118.16	0.51	0.43	0.43	0.32	0.37	0.33	0.25	0.28	0.30	0.22	0.21	0.34
iris	64.90	4.71	0.10	0.10	0.18	0.11	0.11	0.09	0.10	0.09	0.11	0.11	0.07	0.11
letter recognition	–	46933.41	772.81	671.20	566.42	187.95	250.25	463.84	282.25	125.69	251.17	86.73	98.78	192.88
liver-disorders	22.83	35.19	0.15	0.22	0.21	0.10	0.10	0.15	0.09	0.12	0.13	0.09	0.08	0.07
magic	–	38392.73	258.70	252.47	299.13	102.03	108.79	269.24	51.91	66.92	88.94	30.97	48.01	95.76
mammography	60.78	10.58	0.18	0.14	0.32	0.10	0.08	0.24	0.08	0.08	0.22	0.08	0.07	0.22
pendigits	–	14089.59	110.48	74.69	96.22	40.31	42.12	58.73	29.49	18.00	25.39	4.30	4.93	20.95
pim	65.38	125.54	0.51	0.68	0.68	0.33	0.30	0.30	0.19	0.20	0.23	0.13	0.17	0.22
segment	279.67	986.58	7.80	7.69	5.96	2.77	2.64	3.87	1.55	2.32	3.68	0.53	0.89	2.38
shuttle	–	145769.48	1105.78	1076.37	2114.85	569.73	619.69	1102.95	342.76	380.65	902.32	343.93	205.27	989.78
sonar	52.86	125.57	0.70	0.84	0.70	0.40	0.24	0.33	0.32	0.30	0.47	0.35	0.38	0.34
thyroids	86.27	28.67	0.89	1.39	1.58	0.43	0.30	0.91	0.22	0.29	0.74	0.34	0.29	0.61
transfusion	36.14	19.28	0.33	0.26	0.33	0.13	0.14	0.23	0.12	0.11	0.14	0.09	0.11	0.16
vehicle	97.22	358.72	1.04	0.92	1.89	0.44	0.45	0.65	0.33	0.36	0.39	0.27	0.25	0.38
vowel	100.60	231.39	1.20	1.39	1.27	0.49	0.82	0.80	0.31	0.37	0.54	0.29	0.29	0.50
waveform	1437.58	22075.68	131.79	97.71	87.48	67.57	21.67	41.65	19.31	17.88	30.62	5.77	10.50	21.90
wdbc	90.45	264.08	0.66	0.75	1.35	0.48	0.65	0.47	0.18	0.34	0.28	0.18	0.25	0.30
wisc	45.49	34.54	0.45	0.50	0.50	0.13	0.21	0.36	0.15	0.16	0.25	0.12	0.12	0.19
wdbc	26.72	63.19	0.79	0.41	0.52	0.32	0.20	0.22	0.26	0.33	0.22	0.19	0.26	0.15
yeast	153.37	400.10	1.59	2.16	2.15	0.73	0.83	0.80	0.20	0.58	0.58	0.23	0.32	0.59