

The global k-means clustering algorithm

Aristidis Likas

Department of Computer Science
University of Ioannina
Greece

Nikos Vlassis

Computer Science Institute
Faculty of Science
University of Amsterdam
The Netherlands

Jacob J. Verbeek

Computer Science Institute
Faculty of Science
University of Amsterdam
The Netherlands

We present the global k-means algorithm which is an incremental approach to clustering that dynamically adds one cluster center at a time through a deterministic global search procedure consisting of N (with N being the size of the data set) executions of the k-means algorithm from suitable initial positions. We also propose modifications of the method to reduce the computational load without significantly affecting solution quality. The proposed clustering methods are tested on well-known data sets and they compare favorably to the k-means algorithm with random restarts.

Keywords: Clustering; K-means algorithm; Global optimization; k -d trees, Data mining.

**Intelligent
Autonomous
Systems**



Contents

1	Introduction	1
2	The global k-means algorithm	1
3	Speeding-up execution	3
3.1	The fast global k-means algorithm	3
3.2	Initialization with k -d trees	4
4	Experimental results	5
5	Discussion and conclusions	7

Intelligent Autonomous Systems

Computer Science Institute

Faculty of Science

University of Amsterdam

Kruislaan 403, 1098 SJ Amsterdam

The Netherlands

tel: +31 20 525 7461

fax: +31 20 525 7490

<http://www.science.uva.nl/research/ias/>**Corresponding author:**

Aristidis Likas

arly@cs.uoi.gr

1 Introduction

A fundamental problem that frequently arises in a great variety of fields such as pattern recognition, image processing, machine learning and statistics is the clustering problem [1]. In its basic form the clustering problem is defined as the problem of finding groups of data points in a given data set. Each of these groups is called a cluster and can be defined as a region in which the density of objects is locally higher than in other regions.

The simplest form of clustering is partitional clustering which aims at partitioning a given data set into disjoint subsets (clusters) so that specific clustering criteria are optimized. The most widely used criterion is the clustering error criterion which for each point computes its squared distance from the corresponding cluster center and then takes the sum of these distances for all points in the data set. A popular clustering method that minimizes the clustering error is the k-means algorithm. However, the k-means algorithm is a local search procedure and it is well-known that it suffers from the serious drawback that its performance heavily depends on the initial starting conditions [2]. To treat this problem several other techniques have been developed that are based on stochastic global optimization methods (eg. simulated annealing, genetic algorithms). However, it must be noted that these techniques have not gained wide acceptance and in many practical applications the clustering method that is used is the k-means algorithm with multiple restarts [1].

In this work we propose the global k-means clustering algorithm, which constitutes a deterministic effective global clustering algorithm for the minimization of the clustering error that employs the k-means algorithm as a local search procedure. The algorithm proceeds in an incremental way: to solve a clustering problem with M clusters, all intermediate problems with $1, 2, \dots, M-1$ clusters are sequentially solved. The basic idea underlying the proposed method is that an optimal solution for a clustering problem with M clusters can be obtained using a series of local searches (using the k-means algorithm). At each local search the $M-1$ cluster centers are always initially placed at their optimal positions corresponding to the clustering problem with $M-1$ clusters. The remaining M -th cluster center is initially placed at several positions within the data space. Since for $M=1$ the optimal solution is known, we can iteratively apply the above procedure to find optimal solutions for all k -clustering problems $k=1, \dots, M$. In addition to effectiveness, the method is deterministic and does not depend on any initial conditions or empirically adjustable parameters. These are significant advantages over all clustering approaches mentioned above.

In the following section starts with a formal definition of the clustering error and a brief description of the k-means algorithm and then describes the proposed global k-means algorithm. Section 3 describes modifications of the basic method that require less computation at the expense of being slightly less effective. Section 4 provides experimental results and comparisons with the k-means algorithm with multiple restarts. Finally Section 5 provides conclusions and describes directions for future research.

2 The global k-means algorithm

Suppose we are given a data set $X = \{x_1, \dots, x_N\}$, $x_n \in R^d$. The M -clustering problem aims at partitioning this data set into M disjoint subsets (clusters) C_1, \dots, C_M , such that a clustering criterion is optimized. The most widely used clustering criterion is the sum of the squared Euclidean distances between each data point x_i and the centroid m_k (cluster center) of the subset C_k which contains x_i . This criterion is called clustering error and depends on the cluster

centers m_1, \dots, m_M :

$$E(m_1, \dots, m_M) = \sum_{i=1}^N \sum_{k=1}^M I(x_i \in C_k) |x_i - m_k|^2 \quad (1)$$

where $I(X) = 1$ if X is true and 0 otherwise.

The k-means algorithm finds locally optimal solutions with respect to the clustering error. It is a fast iterative algorithm that has been used in many clustering applications. It is a point-based clustering method that starts with the cluster centers initially placed at arbitrary positions and proceeds by moving at each step the cluster centers in order to minimize the clustering error. The main disadvantage of the method lies in its sensitivity to initial positions of the cluster centers. Therefore, in order to obtain near optimal solutions using the k-means algorithm several runs must be scheduled differing in the initial positions of the cluster centers.

In this paper, the *global k-means* clustering algorithm is proposed, which constitutes a deterministic global optimization method that does not depend on any initial parameter values and employs the k-means algorithm as a local search procedure. Instead of randomly selecting initial values *for all cluster centers* as is the case with most global clustering algorithms, the proposed technique proceeds in an incremental way attempting to optimally add one new cluster center at each stage.

More specifically, to solve a clustering problem with M clusters the method proceeds as follows. We start with one cluster ($k = 1$) and find its optimal position which corresponds to the centroid of the data set X . In order to solve the problem with two clusters ($k = 2$) we perform N executions of the k-means algorithm from the following initial positions of the cluster centers: the first cluster center is always placed at the optimal position for the problem with $k = 1$, while the second center at execution n is placed at the position of the data point x_n ($n = 1, \dots, N$). The best solution obtained after the N executions of the k-means algorithm is considered as the solution for the clustering problem with $k = 2$. In general, let $(m_1^*(k), \dots, m_k^*(k))$ denote the final solution for k -clustering problem. Once we have found the solution for the $(k - 1)$ -clustering problem, we try to find the solution of the k -clustering problem as follows: we perform N runs of the k-means algorithm with k clusters where each run n starts from the initial state $(m_1^*(k - 1), \dots, m_{(k-1)}^*(k - 1), x_n)$. The best solution obtained from the N runs is considered as the solution $(m_1^*(k), \dots, m_k^*(k))$ of the k -clustering problem. By proceeding in the above fashion we finally obtain a solution with M clusters having also found solutions for all k -clustering problems with $k < M$.

The latter characteristic can be advantageous in many applications where the aim is also to discover the ‘correct’ number of clusters. To achieve this, one has to solve the k -clustering problem for various numbers of clusters and then employ appropriate criteria for selecting the most suitable value of k [3]. In this case the proposed method directly provides clustering solutions for all intermediate values of k , thus requiring no additional computational effort.

In what concerns computational complexity, the method requires N executions of the k-means algorithm for each value of k ($k = 1, \dots, M$). Depending on the available resources and the values of N and M , the algorithm may be an attractive approach, since, as experimental results indicate, the performance of the method is excellent. Moreover, as we will show later there are several modifications that can be applied in order to reduce the computational load.

The rationale behind the proposed method is based on the following assumption: an optimal clustering solution with k clusters can be obtained through local search (using k-means) starting from an initial state with

- the $k - 1$ centers placed at the optimal positions for the $(k - 1)$ -clustering problem and
- the remaining k -th center placed at an appropriate position to be discovered.

This assumption seems very natural: we expect that the solution of the k -clustering problem to be *reachable* (through local search) from the solution of $(k - 1)$ -clustering problem, once the additional center is placed at an appropriate position within the data set. It is also reasonable to restrict the set of possible initial positions of the k -th center to the set X of available data points. It must be noted that this is a rather computational heavy assumption and several other options (examining fewer initial positions) may also be considered. The above assumptions are also verified experimentally, since in all experiments (and for all values of k) the solution obtained by the proposed method was at least as good as that obtained using numerous random restarts of the k -means algorithm. In this spirit, we can cautiously state that the proposed method is *experimentally optimal* (although it is difficult to prove theoretically).

3 Speeding-up execution

Based on the general idea of the global k -means algorithm, several heuristics can be devised to reduce the computational load without significantly affecting the quality of the solution. In the following subsections two modifications are proposed, each one referring to a different aspect of the method.

3.1 The fast global k -means algorithm

The fast global k -means algorithm constitutes a straightforward method to accelerate the global k -means algorithm. The difference lies in the way a solution for the k -clustering problem is obtained, given the solution of the $(k - 1)$ -clustering problem. For each of the N initial states $(m_1^*(k - 1), \dots, m_{(k-1)}^*(k - 1), x_n)$ we do not execute the k -means algorithm until convergence to obtain the final clustering error E_n . Instead we compute an *upper* bound $E_n \leq E - b_n$ on the resulting error E_n for all possible allocation positions x_n , where E is the error in the $(k - 1)$ -clustering problem. We then initialize the position of the new cluster center at the point x_i that minimizes E_n , or equivalently that maximizes b_n , and execute the k -means algorithm to obtain the solution with k clusters. Formally we have

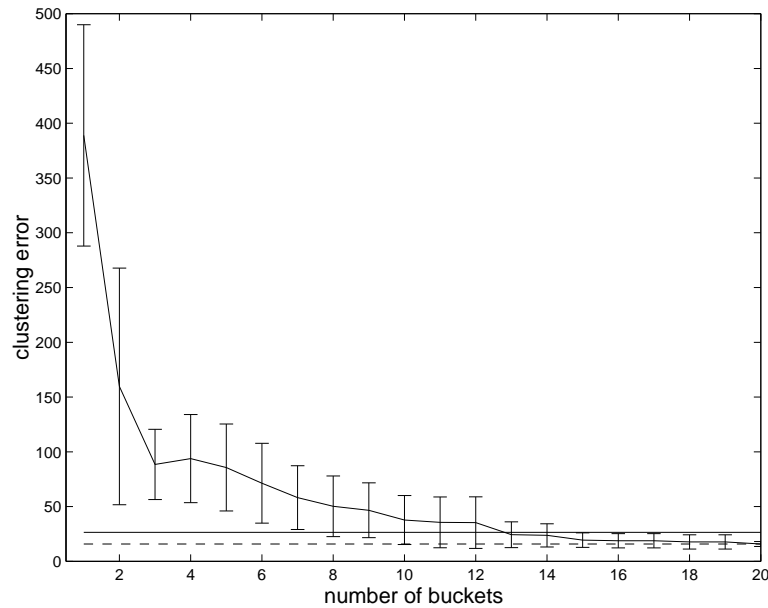
$$b_n = \sum_{j=1}^N \max(d_{k-1}^j - |x_n - x_j|^2, 0), \quad (2)$$

$$i = \arg \max_n b_n \quad (3)$$

where d_{k-1}^j is the squared distance between x_j and the closest center among the $k - 1$ cluster centers obtained so far (ie., center of the cluster where x_j belongs). The quantity b_n measures the *guaranteed* reduction in the error measure obtained by inserting a new cluster center at position x_n .

Suppose the solution of the $(k - 1)$ -clustering problem is $(m_1^*(k - 1), \dots, m_{(k-1)}^*(k - 1))$ and a new cluster center is added at location x_n . Then the new center will allocate all points x_j whose squared distance from x_n is smaller than the distance d_{k-1}^j from their previously closest center. Therefore, for each such data point x_j the clustering error will decrease by $d_{k-1}^j - |x_n - x_j|^2$. The summation over all such data points x_j provides the quantity b_n for a specific insertion location x_n . Since the k -means algorithm is guaranteed to decrease the clustering error at each step, $E - b_n$ upper bounds the error measure that will be obtained if we run the algorithm until convergence after inserting the new center at x_n (this is the error measure used in the global k -means algorithm).

Experimental results (see next section) suggest that using the data point that minimizes this bound leads to results almost as good as those provided by the global k -means algorithm.



Moreover, the cluster insertion procedure can be efficiently implemented by storing in a matrix all pairwise squared distances between points when the algorithm starts, and using this matrix for directly computing the upper bounds above. A similar ‘trick’ has been used in the related problems of greedy mixture density estimation using the EM algorithm [4] and principal curve fitting [5].

Finally, we may still apply this method as well as the global k -means algorithm when we do not consider every data point x_n ($n = 1, \dots, N$) as possible insertion position for the new center, but use only a smaller set of appropriately selected insertion positions. A fast and sensible choice for selecting such a set of positions based on k -d trees is discussed next.

3.2 Initialization with k -d trees

A k -d tree [6, 7] is a multi-dimensional generalization of the standard one-dimensional binary search tree, that facilitates storage and search over k -dimensional data sets. A k -d tree defines a recursive partitioning of the data space into disjoint subspaces. Each node of the tree defines a subspace of the original data space and, consequently, a subset containing the data points residing in this subspace. Each nonterminal node has two successors, each of them associated with one of the two subspaces obtained from the partitioning of the parent space using a cutting hyperplane. The k -d tree structure was originally used for speeding up distance-based search operations like nearest neighbors queries, range queries, etc.

In our case we use a variation of the original k -d tree proposed in [7]. There, the cutting hyperplane is defined as the plane that is perpendicular to the direction of the principal component of the data points corresponding to each node, therefore the algorithm can be regarded as a method for *nested* (recursive) principal component analysis of the data set. The recursion usually terminates if a terminal node (called bucket) is created containing less than a prespecified number of points b (called bucket size) or if a prespecified number of buckets have been created. It turns out that, even if the algorithm is not used for nearest neighbor queries, merely the construction of the tree provides a very good preliminary clustering of the data set. The idea is to use the bucket centers (which are fewer than the data points) as possible insertion locations for the algorithms presented previously.

In Figure 1 average performance results are shown on 10 data sets each one consisting of 300 data points drawn from the same mixture of 15 Gaussian components. The components of the

Gaussian mixture are well separated and exhibit limited eccentricity.

We compare the results of three methods to the clustering problem with $k = 15$ centers: (i) The dashed line depicts the results when using the fast global k-means algorithm with all data points constituting potential insertion locations. The average clustering error over the 10 data sets is 15.7 with standard deviation 1.2. (ii) The solid line depicts results when the standard k-means algorithm is used: one run for each data set was conducted. At each run the 15 cluster centers were initially positioned to the centroids of the buckets obtained from the application of the k -d tree algorithm until 15 buckets were created. The average clustering error over the 10 data sets is 24.4 with standard deviation 9.8. (iii) The solid line (with error bars) shows the results when using the fast global k-means algorithm with the potential insertion locations constrained by the centroids of the buckets of a k -d tree. On the horizontal axis we vary the number of buckets for the k -d tree of the last method.

We also computed the ‘theoretical’ clustering error for each data set, ie., the error computed by using the true cluster centers. The average error value over the 10 data sets was 14.9 with standard deviation 1.3. These results were too close to the results of the standard fast global k-means to include them in the figure.

We can conclude from this experiment that (a) the fast global k-means approach gives rise to performance significantly better than when starting with all centers at the same time initialized using the k -d tree method, and (b) restricting the insertion locations for the fast global k-means to those given by the k -d tree (instead of using all data points) does not significantly degrade performance if we consider a sufficiently large number of buckets in the k -d tree (in general larger than the number clusters).

Obviously, it is also possible to employ the above presented k -d tree approach with the global k-means algorithm.

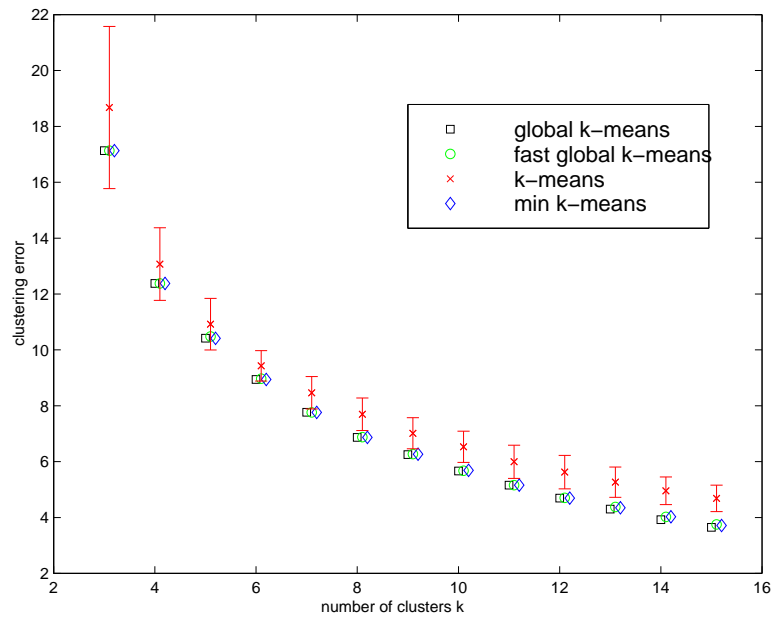
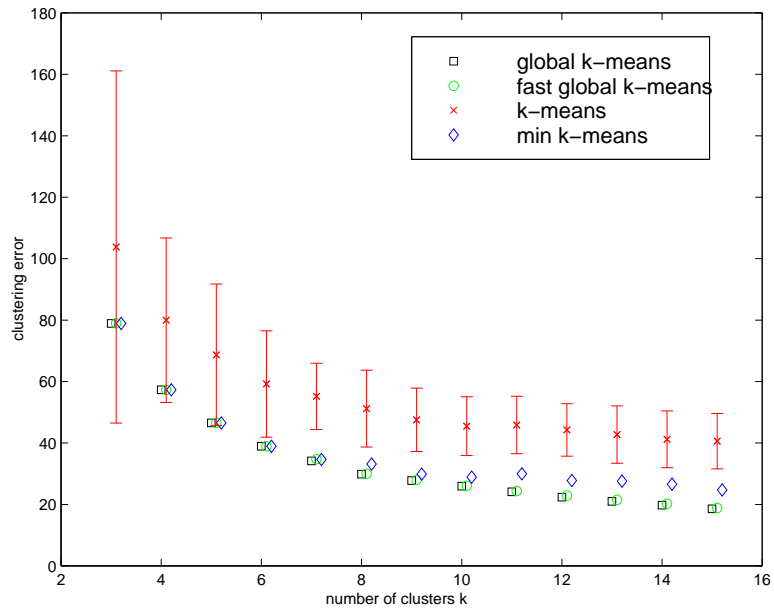
4 Experimental results

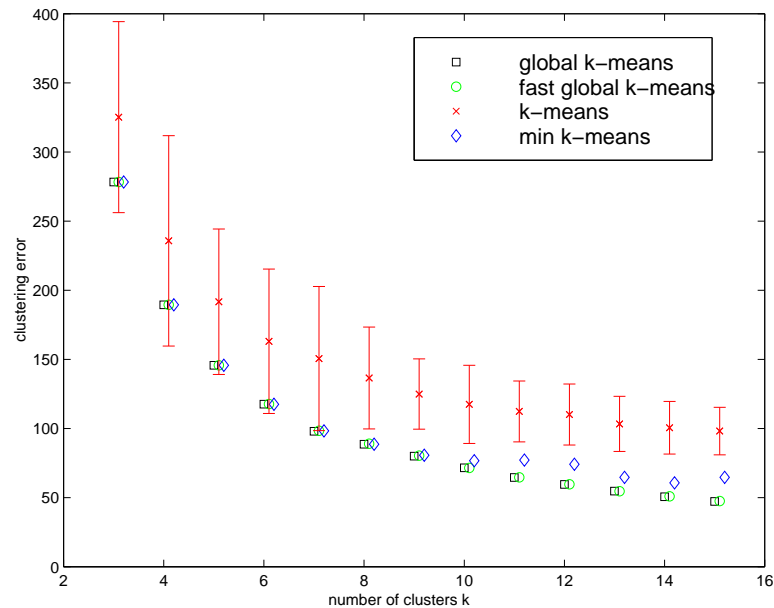
We have tested the proposed clustering algorithms on several well-known data sets, namely the iris data set [8], the synthetic data set [9] and the image segmentation data set [8]. In all data sets we conducted experiments for the clustering problems obtained by considering only feature vectors and ignoring class labels. The iris data set contains 150 four-dimensional data points, the synthetic data set 250 two-dimensional data points and the for the image segmentation data set we consider 210 six-dimensional data points obtained through PCA on the original 18-dimensional data points. The quality of the obtained solutions was evaluated in terms of the values of the final clustering error.

For each data set we conducted the following experiments:

- one run of the global k-means algorithm for $M = 15$.
- one run of the fast global k-means algorithm for $M = 15$.
- the k-means algorithm for $k = 1, \dots, 15$. For each value of k , the k-means algorithm was executed N times (where N is the number of data points) starting from random initial positions for the k centers and we computed the minimum and average clustering error as well as its standard deviation.

For each of the three data sets the experimental results are displayed in Figures 2, 3 and 4 respectively. Each figure plot displays the clustering error value as a function of the number of clusters. It is clear that the global k-means algorithm is very effective providing in all cases solutions of equal or better quality with respect to the k-means algorithm. In what concerns the fast version of the algorithm, it is very encouraging that, although executing significantly





faster, it provides solutions of excellent quality, comparable to those obtained by the original method. Therefore, it constitutes a very efficient algorithm, both in terms of solution quality and computational complexity and can run even faster if k -d trees are employed as explained in the previous section.

Matlab implementations of the fast global k-means and the k -d tree building algorithms can be downloaded from <http://www.science.uva.nl/research/ias>.

5 Discussion and conclusions

We have presented the global k-means clustering algorithm, which constitutes a deterministic clustering method providing excellent results in terms of the clustering error criterion. The method is independent of any starting conditions and compares favorably to the k-means algorithm with multiple random restarts. The deterministic nature of the method is particularly important in cases where the clustering method is used either to specify initial parameter values for other methods (for example RBF training) or constitutes a module in a more complex system. In such a case we can be almost certain that the employment of the global k-means (or any of the fast variants) will always provide sensible clustering solutions. Therefore, one can evaluate the complex system and adjust critical system parameters without having to worry for dependence of system performance on the clustering method employed.

Another advantage of the proposed technique is that in order to solve the M -clustering problem, all intermediate k -clustering problems are also solved for $k = 1, \dots, M$. This may prove useful in many applications where we seek for the actual number of clusters and the k -clustering problem is solved for several values of k . We have also developed the fast global k-means algorithm, which significantly reduces the required computational effort, while at the same time providing solutions of almost the same quality.

We have also proposed two modification of the method that reduce the computational load without significantly affecting solution quality. These methods can be employed to find solutions to clustering problems with thousands of high-dimensional points and one of our primary aims is to test the techniques on large scale data mining problems.

Another direction of future work is related with the use of parallel processing for accelerating the proposed methods since, for every k , the N executions of the k-means algorithm are inde-

pendent and can be performed in parallel. Another research direction concerns the application of the proposed method to other types of clustering (for example fuzzy clustering), as well as to topographic methods like SOM. Moreover, an important issue that deserves further study is related with the possible development of theoretical foundations for the assumptions behind the method. Finally, it is also possible to employ the global k-means algorithm as a method for providing effective initial parameter values for RBF networks and data modeling problems using Gaussian mixture models and compare the effectiveness of the obtained solutions with other training techniques for Gaussian mixture models [10, 11].

Acknowledgements

N. Vlassis and J. J. Verbeek are supported by the Dutch Technology Foundation STW project AIF 4997.

References

- [1] M. N. Murty A. K. Jain and P. J. Flynn, "Data clustering: A review," *ACM Computing Surveys*, vol. 31, no. 3, pp. 264–323, 1999.
- [2] J. A. Lozano J. M. Pena and P. Larranaga, "An empirical comparison of four initialization methods for the k-means algorithm," *Pattern Recognition Letters*, vol. 20, pp. 1027–1040, 1999.
- [3] G. W. Milligan and M. C. Cooper, "An examination of procedures for determining the number of clusters in a data set," *Psychometrika*, vol. 50, pp. 159–179, 1985.
- [4] N. Vlassis and A. Likas, "A greedy EM algorithm for Gaussian mixture learning," Tech. Rep., Computer Science Institute, University of Amsterdam, The Netherlands, Sept. 2000, IAS-UVA-00-08.
- [5] J.J. Verbeek, N. Vlassis, and B. Kröse, "A k -segments algorithm to find principal curves," Tech. Rep., Computer Science Institute, University of Amsterdam, The Netherlands, Nov. 2000, IAS-UVA-00-11.
- [6] J. L. Bentley, "Multidimensional binary search trees used for associative searching," *Commun. ACM*, vol. 18, no. 9, pp. 509–517, 1975.
- [7] R. F. Sproull, "Refinements to nearest-neighbor searching in k-dimensional trees," *Algorithmica*, vol. 6, pp. 579–589, 1991.
- [8] C. L. Blake and C. J. Merz, "UCI repository of machine learning databases," University of California, Irvine, Dept. of Information and Computer Sciences, 1998.
- [9] B. D. Ripley, *Pattern Recognition and Neural Networks*, Cambridge University Press, Cambridge, U.K., 1996.
- [10] G. J. McLachlan and D. Peel, *Finite Mixture Models*, Wiley, New York, 2000.
- [11] N. Vlassis and A. Likas, "A kurtosis-based dynamic approach to Gaussian mixture modeling," *IEEE Trans. Systems, Man, and Cybernetics, Part A*, vol. 29, no. 4, pp. 393–399, July 1999.



***Intelligent
Autonomous
Systems***

IAS reports

This report is in the series of IAS technical reports. The series editor is Stephan ten Hagen (stephanh@science.uva.nl). Within this series the following titles appeared:

See: <http://www.science.uva.nl/research/ias/tr/>

You may order copies of the IAS technical reports from the corresponding author or the series editor. Most of the reports can also be found on the web pages of the IAS group (see the inside front page).



***Intelligent
Autonomous
Systems***