
Clustering in the membership embedding space

Maurizio Filippone

Department of Computer Science,
University of Sheffield,
Regent Court, 211 Portobello Street, Sheffield S1 4DP, UK
E-mail: filippone@dcs.shef.ac.uk

Francesco Masulli*

DISI, Dipartimento di Informatica e Scienze dell' Informazione,
Università di Genova and CNISM,
Via Dodecaneso 35, Genoa, Italy
and
Center for Biotechnology,
Temple University,
1900 N 12th Street, Philadelphia, PA 19122, USA
E-mail: masulli@disi.unige.it
*Corresponding author

Stefano Rovetta

DISI, Dipartimento di Informatica e Scienze dell' Informazione,
Università di Genova and CNISM,
Via Dodecaneso 35, Genoa, Italy
E-mail: rovetta@disi.unige.it

Abstract: In several applications of data mining to high-dimensional data, clustering techniques developed for low-to-moderate sized problems obtain unsatisfactory results. This is an aspect of the *curse of dimensionality* issue. A traditional approach is based on representing the data in a suitable similarity space instead of the original high-dimensional attribute space. In this paper, we propose a solution to this problem using the projection of data onto a so-called membership embedding space obtained by using the memberships of data points on fuzzy sets centred on some prototypes. This approach can increase the efficiency of the popular fuzzy C-means method in the presence of high-dimensional datasets, as we show in an experimental comparison. We also present a constructive method for prototypes selection based on simulated annealing that is viable for semi-supervised clustering problems.

Keywords: high-dimensional datasets; unsupervised clustering; semi-supervised clustering; fuzzy sets; embedding spaces; fuzzy C-means; FCM; simulated annealing; SA; curse of dimensionality; knowledge engineering.

Reference to this paper should be made as follows: Filippone, M., Masulli, F. and Rovetta, S. (2009) 'Clustering in the membership embedding space', *Int. J. Knowledge Engineering and Soft Data Paradigms*, Vol. 1, No. 4, pp.363–375.

Biographical notes: Maurizio Filippone received his Master's in Physics in 2004 and his PhD in Computer Science in 2008 from the University of Genova. In 2007, he has been a Research Scholar at the Information and Software Engineering Department at George Mason University. In 2008, he joined the Machine Learning group at the University of Sheffield as a Research Associate. His research interests are focused on machine learning techniques and their applications to pattern recognition, bioinformatics and time series analysis and forecasting.

Francesco Masulli received his Laurea in Physics in 1976 from the University of Genova. He is currently an Associate Professor in Computer Science at the University of Genova. From 1983 to 2001, he has been an Assistant Professor at the University of Genova and from 2001 to 2005, an Associate Professor at University of Pisa. He authored or co-authored more than 120 scientific papers in machine learning, neural networks, fuzzy systems, pattern recognition and bioinformatics.

Stefano Rovetta received his Laurea in Electronic Engineering in 1992 and his PhD in Models, Methods and Tools for Electronic and Electromagnetic Systems in 1996, both from the University of Genova. He is currently an Assistant Professor in Computer Science at the University of Genova. His current research interests include machine learning and clustering techniques for high-dimensional biomedical data analysis and document analysis.

1 Introduction

Clustering methods are useful tools for data mining. They can be employed both in an *unsupervised* way, when available data are unlabelled (or available labels are unreliable, or when the data labelling task is too expensive) and in a *semi-supervised* way when a small amount of knowledge is available concerning either pairwise (must-link or cannot-link) constraints between data items or class labels for some items.

Unfortunately, in several problems of data mining, data lie in a very high-dimensional space (thousands of dimensions). In these situations, the direct application of clustering algorithms developed for low-dimensional spaces [e.g., K-means (Steinhaus, 1956; Lloyd, 1982), BIRCH (Zhang et al., 1996), CURE (Guha et al., 1998), CLARANS (Ng and Han, 2002), etc.] often leads to poor results. Even after feature selection, one may be left with hundreds of dimensions (and further reductions will significantly degrade the results). Even techniques designed for large quantities of data, such as DBScan (Ester et al., 1996), only focus on the problem of large-cardinality datasets, thus effectively making the assumption of (relatively) small dimensionality. This is an aspect of the well-known *curse of dimensionality* issue (Bellman, 1961).

Many clustering algorithms suffer from being applied in high-dimensional spaces, as clustering algorithms often seek for areas where data are dense. Sometimes the cardinality of the datasets available is even less than the number of variables, as in the case of the analysis of many bioinformatics datasets or in web mining problems. This means that data span only a subspace within the data space. In these conditions, it is not easy to define the concept of volumetric density.

Moreover, when space dimensionality is high or even moderate (as low as 10–15), the distance of a point to its farthest neighbour and to its nearest neighbour tend to become

equal (Beyer et al., 1999; Aggarwal and Yu, 2002). Therefore, the evaluation of distances and the concept of *nearest neighbour* itself, become less and less meaningful with the growing dimensions. Defining clusters on the basis of distance measures requires that distances can be estimated. For instance, one of the most commonly used methods, K-means clustering (Steinhaus, 1956; Lloyd, 1982), is based on iteratively computing distances and cluster averages. Increasing the data space dimensionality may introduce a large number of suboptimal solutions (local minima) and the nearest-neighbour criterion which is the basis of the method may even become useless. This problem is not avoided even when K-means is modified in the direction of incorporating fuzzy concepts, e.g., as for the fuzzy C-means (FCM) algorithm (Dunn, 1973; Bezdek, 1981).

A possible approach alleviating these problems is based on representing the data in a suitable similarity space instead of the original high-dimensional attribute space (see, e.g., Strehl and Ghosh, 2003; Filippone et al., 2008; Filippone, 2009).

In this paper, we propose a solution to the highlighted problems using the projection of data onto a so-called membership embedding space (MES). Such projection is obtained by using the memberships of data points on fuzzy sets centred on some prototypes selected among data points themselves. We will demonstrate that this approach can increase clustering efficiency of the popular FCM (Bezdek, 1981) algorithm in the presence of high-dimensional datasets. To this aim, we will experimentally compare the performances of FCM in the original data space, with those in the distance embedding space (DES) [following the approach proposed by Pekalska et al. (2001)] and MES, using different prototypes-data ratios. Moreover, we will present a constructive method for prototypes selection based on simulated annealing (SA) that is viable for semi-supervised clustering problems as well.

In Section 2, a fuzzy embedding for high-dimensional datasets is presented; in Section 3, we recall the main aspects of FCM clustering algorithm; and in Section 4, we present a constructive approach for selecting an optimal set of prototypes in the fuzzy embedding. The experimental results are reported in Section 5. The conclusions are given in Section 6.

2 Membership embedding space

A notable complexity reduction of data mining problems in the presence of large dimensional datasets can be provided by representations in a similarity space or embedding space based on an assigned pairwise similarity (or dissimilarity) transformation (see, e.g., Strehl and Ghosh, 2003; Filippone et al., 2008; Filippone, 2009).

Given a dataset X of cardinality n , $X = \{x_1, x_2, \dots, x_n\}$ in a d -dimensional space, the (dis-)similarity transformation $v(x_i, x_j)$ maps the $n \times d$ data matrix into a more dense symmetric $n \times n$ matrix of similarities v , with $v_{ik} = v(x_i, x_k) \forall i, k$.

Mutual distances or other pairwise pattern evaluation methods such as kernels (Shawe-Taylor and Cristianini, 2004) may be used as (dis-)similarity transformations. If the cardinality of the dataset is small compared to the input space dimensionality, datasets can be represented in the embedding space in a very compact way.

Applications of projection onto (dis-)similarity embedding spaces to clustering are reported, e.g., in Fred and Leitão (2003) and Rovetta and Masulli (2006). Pekalska et al.

(2001) developed a set of methods based on representing each pattern according to a set of similarity measurements with respect to other patterns in the dataset. As they pointed out, the (dis-)similarity measure should be a metric, since metrics preserve the *reverse of the compactness hypothesis* (Pekalska et al., 2001): ‘objects that are similar in their representation are also similar in reality and belong, thereby, to the same class’.

Often, non-metric distances are used as well. Moreover, sometimes the (dis-)similarity matrix can be reduced from a square matrix $n \times n$ to a smaller rectangular matrix $n \times s$, by selecting $s \leq n$ reference points (called *prototypes*) and computing the (dis-)similarities of the data with them. If the embedding dimension s is small compared with d (i.e., $s/d \ll 1$), some points could have an ambiguous representation.

In order to avoid the previously highlighted problems, in this paper, we study an embedding based on the space of memberships to fuzzy sets (Zadeh, 1965) centred on selected prototypes.

The memberships to fuzzy sets centred on the prototypes are modelled using the following normalised function:

$$v_{ik} = \frac{\exp[-\beta d_{ik}^2]}{\sum_l \exp[-\beta d_{lk}^2]} \quad \beta \in \mathfrak{R}^+, \quad (1)$$

where the parameter β regulates the spread of the membership function. The denominator normalises the sum of the memberships to the prototypes to sum up to one. The matrix $\mathbf{V} = [v_{ik}]$ is the similarity matrix. Note that \mathbf{V} is rectangular since we select a number of prototypes $s \leq n$.

In this way, in the MES a data point x_i is represented as a row of v , i.e., $x_i = (v_{i1}, v_{i2}, \dots, v_{im})$. Due to the localised definition of fuzzy sets, this vector of memberships contains only few non-null elements in correspondence of the nearest prototypes in the original data space. If the spread of membership is large (i.e., large β) many of these elements are non-null, otherwise for β going to zero, only the data points corresponding to the selected prototypes have at least one non-null element.

The results of a clustering method will be affected by the number and the positions of the prototypes as well as by the value of spread parameter β . Placing the prototypes is a combinatorial search problem which will be tackled by an SA approach.

3 FCM algorithm

In the experiments reported in this paper, we have used the FCM algorithm (Bezdek, 1981) as the clustering algorithm. Other clustering techniques can be applied, but we focus on a single choice for the sake of clarity.

The FCM algorithm performs the minimisation of the following functional:

$$J_m(\mathbf{U}, \mathbf{Y}) \equiv \sum_{i=1}^n \sum_{k=1}^c (u_{ik})^m d_{ik} \quad (2)$$

where $X = \{x_1, x_2, \dots, x_n\}$ is a dataset containing n unlabelled sample points; $Y = \{y_1, y_2, \dots, y_c\}$ is the set of the centres of clusters; $\mathbf{U} = [u_{ik}]$ is the $c \times n$ fuzzy c -partition matrix, containing the membership values of all samples to all prototypes;

$m \in (1, \infty)$ is the fuzziness control parameter; and d_{ik} is a dissimilarity measure between data point x_i and the centre y_k of a specific cluster k . In the rest of this paper, we will use the Euclidean squared distance $d_{ik} \equiv \|x_i - y_k\|^2$ as the dissimilarity measure.

The clustering problem can be formulated as the minimisation of J_m with respect to Y , under the normalisation constraint $\sum_{k=1}^c u_{ik} = 1$.

Then the necessary conditions for minimising J_m are:

$$y_k = \frac{\sum_{i=1}^n (u_{ik})^m x_i}{\sum_{i=1}^n (u_{ik})^m} \quad \text{for all } k, \quad (3)$$

$$u_{ik} = \left[\sum_{j=1}^c \left(\frac{d_{ik}}{d_{jk}} \right)^{\frac{1}{m-1}} \right]^{-1} \quad \text{for all } i, k. \quad (4)$$

The FCM algorithm usually starts with a random initialisation of the fuzzy c-partition matrix \mathbf{U} or of the centroids y_k . Then, it iterates equations (3) and (4) until convergence. Usually, the convergence is checked by comparing the change in the position of the centroids or in the cost function with some fixed thresholds.

Note that in the limit for $m \rightarrow 1$ the FCM functional J_m [equation (2)] tends to the expectation of the K-means global error $\langle E \rangle \equiv \sum_{i=1}^n \sum_{k=1}^c u_{ik} d_{ik}$ and the FCM behaves as the classic K-means (or hard C-means) algorithm (Steinhaus, 1956; Lloyd, 1982; Duda and Hart, 1973).

4 Simulated annealing prototype selection algorithm

As already noted, the selection of the optimal set of prototypes for constructing the MES is a combinatorial search problem. A constructive heuristic algorithm able to select the set of prototypes leading to (sub-)optimal clustering in the MES can be based on SA (Kirkpatrick et al., 1983; Černý, 1985) that is a global search probabilistic technique inspired to annealing in metallurgy.

‘Physical’ annealing as used in metallurgy involves heating a material and then cooling it slowly and in a controlled fashion. The aim of this process is to allow the crystal lattice to reorganise so as to reduce the defects and to reach a more stable and, therefore, stronger inner structure. Heating allows atoms to detach from their initial positions (corresponding to a local minimum of the internal energy) and to float randomly through states of higher energy, slow cooling allows them more chances to find configurations with internal energy lower than the initial one.

SA is an adaptation of the Metropolis-Hastings algorithm (Metropolis et al., 1953) aimed to simulate the behaviour and small fluctuations of a system of atoms starting from an initial configuration, by the generation of a sequence of iterations. In the Metropolis algorithm, each iteration comprises a random perturbation (modification) of the actual configuration (state vector) and the computation of the corresponding energy variation

(ΔE). If $\Delta E < 0$, the transition is unconditionally accepted, otherwise the transition is accepted with probability given by the Boltzmann distribution:

$$P(\Delta E) = \exp\left(\frac{-\Delta E}{KT}\right) \quad (5)$$

where K is the Boltzmann constant and T the temperature.

In SA, this approach is generalised to the solution of general optimisation problems (Kirkpatrick et al., 1983) by using an *ad hoc* selected cost function (*generalised energy*) instead of the physical energy, therefore, it can also be employed when the search space is discrete, as in combinatorial search problems. SA works as a probabilistic hill-climbing procedure searching for the global optimum of the cost function (Romeo, 1986). K is usually set to 1, while the temperature T controls the size of the search area and is gradually lowered until no further improvements of the cost function are noticed. SA can work in very high-dimensional searches, given enough computational resources. In applications, it is important to trade-off the quality of the solution and the computational cost; a slower decreasing of the temperature allows the system to reach better solutions but more time is required to explore the state space.

In Table 1 the proposed simulated annealing prototype selection (SA-PS) algorithm is shown. In our approach, the state of the system is represented by a binary mask $\mathbf{g} = (g_1, g_2, \dots, g_n)$, where each bit g_i (with $i = 1, \dots, n$) corresponds to the selection ($g_i = 1$)/deselection ($g_i = 0$) of a prototype. The initialisation of the vector mask \mathbf{g} (step 2) is done by generating s_0 integer numbers with uniform distribution in the interval $[1, n]$ and setting the corresponding bits of \mathbf{g} to 1 and the remaining ones to 0. At each step, only s prototypes are selected from the original set of n patterns. A perturbation or move is done in the following way:

- 1 choose randomly $w \in [w_{\min}, w_{\max}]$ and $v \in [v_{\min}, v_{\max}]$
- 2 w bits of \mathbf{g} set to 1 are switched to 0
- 3 v bits of \mathbf{g} set to 0 are switched to 1.

The values w_{\min} , w_{\max} , v_{\min} , v_{\max} can be used to reduce or to increase the variability of each perturbation.

Once a set of prototypes is selected, it is possible to represent each pattern in the MES and perform clustering.

The generalised energy E is computed as a linear combination between an assigned clustering quality measure ε and the number of selected prototypes s :

$$E = \varepsilon + \lambda s \quad (6)$$

The clustering quality measure ε can be a function of either the cost function associated to the clustering algorithm, a clustering validation index or, in the case of *semi-supervised* clustering (where we have a partially labelled dataset), the *representation error* (RE). RE is the number of data points in each cluster disagreeing with the majority label in that cluster, summed over all clusters and expressed as a percentage.

Note that the introduction of the number of selected prototypes s in the computation of E penalises situations in which the number of selected prototypes is high, effectively resulting as a complexity penalty term. This choice of E leads to the minimisation of the

cardinality of the set of prototypes able to achieve a good clustering quality measure. The balance between these two terms is controlled by λ (*penalisation coefficient*).

The cooling strategy is implemented in step 6 of Table 1. It should be noted that this strategy is only one of the many possible choices. The application of different strategies and in fact, the value of α itself, can significantly influence the quality of results and the computing time. The decay law and the parameter value used have proven to be reasonable in our experiments, but they may have to be evaluated on the specific application.

Table 1 SA-PS algorithm

1	Initialise parameters (see list in Table 3);
2	Initialise the binary mask \mathbf{g} at random;
3	Perform clustering and evaluate the generalised system energy E ;
4	do
5	Initialise $f=0$ (number of iterations, $h=0$ (number of successes);
a	do
b	Increment number of iterations f ;
c	Pertube mask \mathbf{g} ;
d	Perform clustering and evaluate the generalised system energy E ;
e	Generate a random number rnd in the interval $[0, 1]$
f	if $rnd < P(\Delta E)$ then
1	Accept the new \mathbf{g} mask
2	Increment the number of successes h ;
g	end if
h	loop until $h \leq h_{\min}$ and $f \leq f_{\max}$;
6	update $T = \alpha T$;
7	loop until $h > 0$
8	end.

5 Experimental results

5.1 Dataset

In order to test our approach, we have used a high-dimensional bioinformatics dataset, the publicly available leukaemia data by Golub et al. (1999). The leukaemia problem consists in characterising two forms of acute leukaemia, acute lymphoblastic leukaemia (ALL) and acute myeloid leukaemia (AML). The original work proposed both a supervised classification task ('class prediction') and an unsupervised characterisation task ('class discovery'). Here we obviously focus on the latter, but we exploit the diagnostic information on the type of leukaemia to assess the goodness of the clustering obtained.

The dataset contains 38 samples for which the expression level of 7129 genes has been measured with the DNA microarray technique (the interesting human genes are

6,817 and the other are controls required by the technique). These expression levels have been scaled by a factor of 100. Of these samples, 27 are cases of ALL and 11 are cases of AML. Moreover, it is known that the ALL class is in fact composed by two different diseases since they are originated from different cell lineages (either T-lineage or B-lineage).

5.2 Performance comparison

We have compared the following approaches:

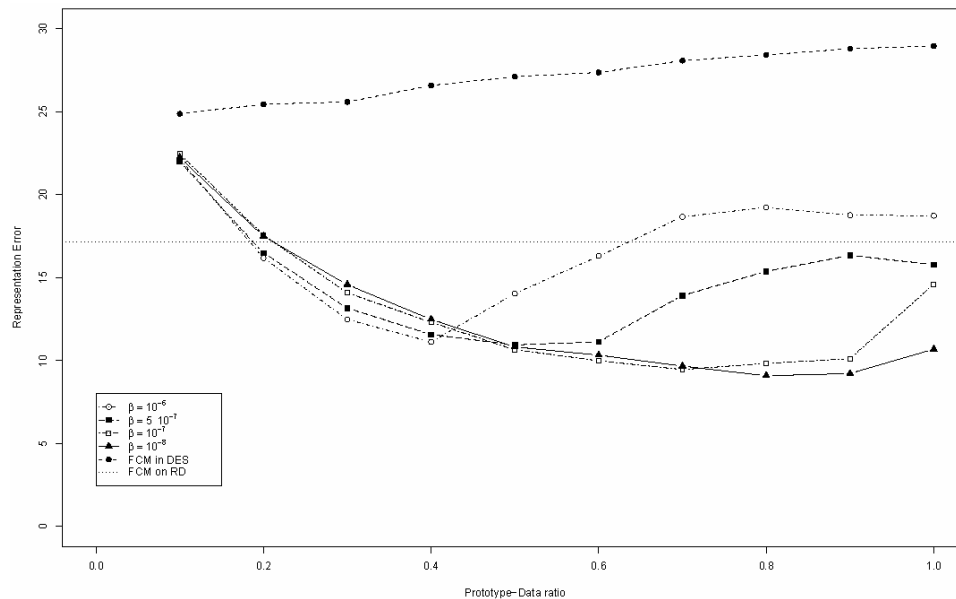
- 1 FCM on the raw dataset (RD)
- 2 FCM in the DES with different prototypes/data ratios
- 3 FCM in the MES with different prototypes/data ratios.

Each experiment corresponds to 1,000 independent trials, each of them using a different random initialisation of memberships in the FCM algorithm. In all trials, the number of clusters was set to 3 and the fuzziness parameter m of FCM was set to 2.

Figure 1 shows the RE versus the prototypes/data ratio averaged over 1,000 independent trials.

The first approach (standard FCM on original data) obtains a RE of 17.2%.

Figure 1 RE for the tested methods



Note: FCM on RD, FCM on the DES and FCM on the MES using $\beta = 10^{-6}$, $5 \cdot 10^{-7}$, 10^{-7} and 10^{-8} .

The projection onto the DES (second approach) leads to worse results compared to the first approach: as we can see for Figure 1, in this case, the RE is greater than 25.0% for all prototypes/data ratios in the range $[.1, 1.0]$.

The last approach, projecting the dataset onto the MES, leads to better results. In Figure 1, we show the results with β starting from 10^{-6} (that is the about the reciprocal of the mean distance between data points) and with decreasing values of this parameter until 10^{-8} that gives the optimal RE.

A comparison of the best RE for the tested methods is reported in Table 2. For each value of β we can notice an optimal prototype/data ratio.

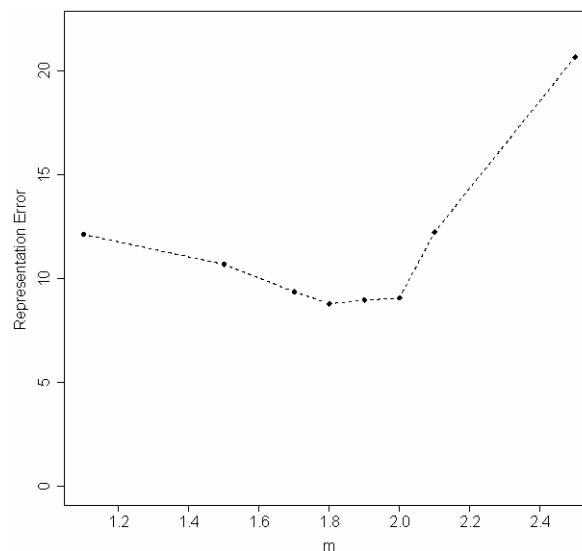
Table 2 Comparison of the best RE for the tested methods

<i>Method</i>	β	<i>RE</i>	<i>Prototypes/data ratio</i>
RD	-	17.2	/
DES	-	24.9	0.1
MES	10^{-6}	11.1	0.4
MES	$5 \cdot 10^{-7}$	10.9	0.5
MES	10^{-7}	9.5	0.7
MES	10^{-8}	9.1	0.8

Note: FCM on RD, FCM on the DES and FCM on the MES.

Finally, we performed a ‘model selection’ to find an appropriate value of the fuzziness m by computing the RE over 1,000 trials. This set of experiments gives also indications about the role of the fuzziness parameter of FCM when applied in an MES. We obtained the MES using $\beta = 10^{-8}$ and prototypes/data ratio = 0.8. As shown in Figure 2, the best value for m is $m = 1.8$ that allows to obtain a RE equal to 8.8% (even if this is slightly better than the results obtained by $m = 2$). For $m > 2$ we notice a rapid increasing in the RE. On the other hand, for low values of m , FCM tends to behave like the K-means algorithm that performs worse than FCM.

Figure 2 The behaviour of the best error rate vs. the fuzziness parameter m



Note: Achieved with $\beta = 10^{-8}$, prototypes/data ratio = 0.8.

5.3 Experiments on the constructive approach

We show here the application of the SA-PS algorithm to the leukaemia data by Golub et al. (1999) simulating a semi-supervised clustering setting.

We ran the SA-PS algorithm in the MES using the FCM (Bezdek, 1981) algorithm to cluster data. As a clustering quality measure, we used the RE evaluated as the best value obtained on $r = 10$ independent trials of FCM.

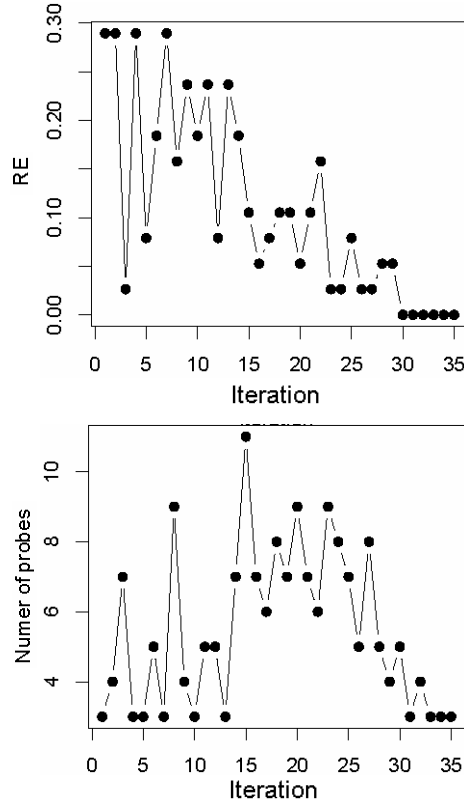
The parameter λ controls the trade-off between the RE and the number of selected probes (that is a measure of complexity). In our case, the penalisation score for each probe corresponds to an RE of 1% ($\lambda = 10^{-2}$). The parameters controlling the annealing are α , f_{\max} and h_{\min} , we selected $\alpha = 0.9$ to allow a slow cooling of the system and $f_{\max} = 2,000$ and $h_{\min} = 200$ in order to have the chance to explore several states for a specific value of T . The number of bits to be switched in each move (w_{\min} , v_{\min} , w_{\max} , v_{\max}) were selected in order to give the system enough variability to perform small as well as long jumps between states.

Table 3 shows the list of parameters of our algorithm and the values we have used in the experiments here reported.

Table 3 SA-SP algorithm – choice of parameters

<i>Meaning</i>	<i>Symbol</i>	<i>Value</i>
Number of random perturbation of \mathbf{g} used to estimate the initial value of T	p	10,000
Number of prototypes to be initially selected	s_0	3
Cooling parameter	α	0.9
Membership width parameter	β	10^{-6}
Maximum number of iteration at each T	f_{\max}	2,000
Minimum number of successes for each T	h_{\min}	200
Penalisation coefficient	λ	10^{-2}
Minimum number of bits to be switched	w_{\min}, v_{\min}	1, 1
Maximum number of bits to be switched	w_{\max}, v_{\max}	$s, 5$
Number of clusters	c	3
FCM fuzziness parameter	m	2
FCM trials	r	10

Each independent run of the SA-PS algorithm finds a different small subset of prototypes leading to a clustering RE equal to zero. In Figure 3, the RE and the number of selected bits of \mathbf{g} are plotted versus the iteration number during a run of the SA-PS algorithm, where each iteration corresponds to a different value of temperature T . In this case, at iterations 31, 33, 34 and 35 we obtained four different sets of three prototypes giving clustering RE equal to zero.

Figure 3 RE and number of prototypes selected during a run of the SA-PS algorithm

6 Conclusions

Clustering methods can achieve poor results when applied to small cardinality and high dimensionality datasets.

In this paper, we proposed a method to face those clustering problems using an embedding space where each data point is represented by a vector containing memberships to fuzzy sets centred on a subset of prototypes selected from the database. On the leukaemia data by Golub et al. (1999) the proposed approach leads to significant improvements with respect the application of clustering algorithms in the original space and in the DES.

The method can exploit supervised information (class labels) even when these are not available for all data points. This is because they are not used in the optimisation step but only in the centroid evaluation step, which is configured as a model selection over centroid position (a 'fitting' criterion) and number (a 'complexity' criterion). This makes the proposed approach a viable solution in all cases where supervised information is available, even if only for a subset of data points.

Obtaining (good quality) supervised information has always been an expensive step in setting up an application, but recently this has become an even more serious issue,

given the enormous quantities of data that can be produced at a fast pace by sources such as, for instance, enterprise data warehouses, the web or high throughput biomolecular analysis techniques. Being able to exploit unsupervised data is important, but perhaps even more important is to be able to exploit even incomplete – but precious – supervised information.

Acknowledgements

This work was funded by the MIUR grant code 2004062740.

References

- Aggarwal, C.C. and Yu, P.S. (2002) ‘Redefining clustering for high-dimensional applications’, *IEEE Transactions on Knowledge and Data Engineering*, Vol. 14, No. 2, pp.210–225.
- Bellman, R. (1961) *Adaptive Control Processes: A Guided Tour*, Princeton University Press.
- Beyer, K.S., Goldstein, J., Ramakrishnan, R. and Shaft, U. (1999) ‘When is ‘nearest neighbor’ meaningful?’, in C. Beeri and P. Buneman (Eds.): *ICDT, Lecture Notes in Computer Science*, Vol. 1540, pp.217–235, Springer.
- Bezdek, J.C. (1981) *Pattern Recognition with Fuzzy Objective Function Algorithms*, Kluwer Academic Publishers, Norwell, MA, USA.
- Černý, V. (1985) ‘Thermodynamical approach to the traveling salesman problem: an efficient simulation algorithm’, *Journal of Optimization Theory and Applications*, Vol. 45, No. 1, pp.41–51.
- Duda, R.O. and Hart, P.E. (1973) *Pattern Classification and Scene Analysis*, John Wiley and Sons.
- Dunn, J.C. (1973) ‘A fuzzy relative of the isodata process and its use in detecting compact well-separated clusters’, *Journal of Cybernetics*, Vol. 3, pp.32–57.
- Ester, M., Kriegel, H-P., Sander, J. and Xu, X. (1996) ‘A density-based algorithm for discovering clusters in large spatial databases with noise’, in *KDD*, pp.226–231.
- Filippone, M. (2009) ‘Dealing with non-metric dissimilarities in fuzzy central clustering algorithms’, *International Journal of Approximate Reasoning*, Vol. 50, No. 2, pp.363–384.
- Filippone, M., Camastra, F., Masulli, F. and Rovetta, S. (2008) ‘A survey of kernel and spectral methods for clustering’, *Pattern Recognition*, Vol. 41, No. 1, pp.176–190.
- Fred, A.L.N. and Leitão, J.M.N. (2003) ‘A new cluster isolation criterion based on dissimilarity increments’, *IEEE Transactions on Pattern Analysis and Machine Intelligence*, Vol. 25, No. 8, pp.944–958.
- Golub, T.R., Slonim, D.K., Tamayo, P., Huard, C., Gaasenbeek, M., Mesirov, J.P., Coller, H., Loh, M.L., Downing, J.R., Caligiuri, M.A., Bloomfield, C.D. and Lander, E.S. (1999) ‘Molecular classification of cancer: class discovery and class prediction by gene expression monitoring’, *Science*, Vol. 286, No. 5439, pp.531–537.
- Guha, S., Rastogi, R. and Shim, K. (1998) ‘Cure: an efficient clustering algorithm for large databases’, in L.M. Haas and A. Tiwary (Eds.): *SIGMOD Conference*, pp.73–84, ACM Press.
- Kirkpatrick, S.D.G. Jr. and Vecchi, M.P. (1983) ‘Optimization by simulated annealing’, *Science*, Vol. 220, No. 4598, pp.671–680.
- Lloyd, S.P. (1982) ‘Least squares quantization in PCM’, *IEEE Transactions on Information Theory*, Vol. IT-28, No. 2, pp.129–137.
- Metropolis, N., Rosenbluth, A.W., Rosenbluth, M.N., Teller, A.H. and Teller, E. (1953) ‘Equation of state calculations by fast computing machines’, *The Journal of Chemical Physics*, Vol. 21, No. 6, pp.1087–1092.

- Ng, R.T. and Han, J. (2002) 'Clarans: a method for clustering objects for spatial data mining', *IEEE Transactions on Knowledge and Data Engineering*, Vol. 14, No. 5, pp.1003–1016.
- Pekalska, E., Paclík, P. and Duin, R.P.W. (2001) 'A generalized kernel approach to dissimilarity-based classification', *Journal of Machine Learning Research*, Vol. 2, pp.175–211.
- Romeo, F. (1986) 'Probabilistic hill climbing algorithms: properties and applications', *Technical Report*, No. UCB/ERL M86/97, EECS Department, University of California, Berkeley.
- Rovetta, S. and Masulli, F. (2006), 'Shared farthest neighbor approach to clustering of high dimensionality, low cardinality data', *Pattern Recognition*, Vol. 39, No. 12, pp.2415–2425.
- Shawe-Taylor, J. and Cristianini, N. (2004) *Kernel Methods for Pattern Analysis*, Cambridge University Press.
- Steinhaus, H. (1956) 'Sur la division des corp materiels en parties', *Bull. Acad. Polon. Sci.*, Vol. 1, pp.801–804.
- Strehl, A. and Ghosh, J. (2003) 'Relationship-based clustering and visualization for high-dimensional data mining', *INFORMS Journal on Computing*, Vol. 15, No. 2, pp.208–230.
- Zadeh, L.A. (1965) 'Fuzzy sets', *Information and Control*, Vol. 8, pp.338–353.
- Zhang, T., Ramakrishnan, R. and Livny, M. (1996) 'Birch: an efficient data clustering method for very large databases', in H.V. Jagadish and I.S. Mumick (Eds.): *SIGMOD Conference*, pp.103–114, ACM Press.