

the corresponding u_{ij} 's are smaller for these vectors, as the updating equation for u_{ij} 's indicates. Nevertheless, while the update of u_{ij} 's is still in the early stages during which their estimates are not good enough yet, the estimates of θ_j 's in each iteration may not be accurate enough at these stages, affecting thus their final estimation. In addition, it can be seen that u_{ij} 's update is highly dependent from the η_j parameter.

3 THE ADAPTIVE PCM (APCM)

In this section, we will describe the modifications that need to be imposed on the PCM algorithm in order to implement the adaptation of both η 's and the number of clusters. Note that the number of clusters m has to be set to a value greater - but not much greater - than the true number of the actual clusters formed in X . Potentially, the algorithm reduces m to the true number of clusters, by leading the cluster representatives in dense regions. Note, also, that even if the number of clusters m is initially underestimated, APCM will manage to uncover at least some of the dense regions, formed naturally in X .

3.1 Parameter initialization

As it has already been mentioned, the initialization of θ_j and η_j is a crucial part of the algorithm, as it dramatically affects the final clustering result. In the proposed clustering scheme, the initialization of θ_j is carried out using the Max-Min algorithm proposed in [9]. Assuming that the data vectors form cohesive clusters, this algorithm will, in principle, return at least one vector from each cluster ([9]), provided that m is sufficiently large. Max-Min algorithm determines m points of X that will be used as initial estimates of the representatives, by first determining the two most distant points of X . Then the Max-Min algorithm proceeds by determining the remaining $m - 2$ data points, as described in the sequel and assigning them to X_{re} , which is the set containing the points that will be used as initial cluster representatives. It determines for each $\mathbf{x} \in X - X_{re}$ its minimum distance $d_m(\mathbf{x})$ from all the points in X_{re} . Then, it identifies the point in $X - X_{re}$ which has the maximum value among all $d_m(\mathbf{x})$ $\mathbf{x} \in X - X_{re}$ and assigns it to X_{re} . The algorithm terminates when X_{re} contains m vectors.

After the initialization of the θ_j 's, we propose η_j 's to be initialized as follows (see also **Algorithm 2**). The distance of each $\theta_j \in X_{re}$ from its closest $\theta_s \in X_{re} - \{\theta_j\}$, denoted by $d_m(\theta_j)$, is determined first and then the algorithm sets $\eta_j = \frac{d_m(\theta_j)/2}{-\log \beta}$, where $\beta \in (0, 1)$ is an appropriately chosen parameter. To unravel the rationale of the last equation, we solve it with respect to β , obtaining

$$\beta = \exp\left(-\frac{d_m(\theta_j)/2}{\eta_j}\right).$$

Algorithm 2 Initialization of η_j 's

```

for j = 1 to m
    Determine  $d_m(\theta_j) = \min_{\theta_s \in X_{re} - \{\theta_j\}} \|\theta_j - \theta_s\|^2$ 
    Set  $\eta_j = \frac{d_m(\theta_j)/2}{-\log \beta}$ 
end for

```

Comparing this with the updating equation for u_{ij} 's of the PCM algorithm, one can deduce that η_j is chosen such that, for a given data point \mathbf{x}_i that lies at distance $d_m(\theta_j)/2$ from θ_j , its degree of compatibility (u_{ij}) with the j th cluster equals to β . As it has been verified experimentally, typical values for β that lead to satisfactory results are in the range 0.5-0.9. The experiments showed also that β depends on how densely the natural clusters are located; low density requires smaller values of β , while high density requires larger values of β .

3.2 Parameter adaptation

Parameter adaptation in the proposed APCM clustering algorithm refers to the adjustment of the number of clusters and the adaptation of η 's, which are two interrelated processes. These processes can be attached in the main while-loop of the classic possibilistic algorithm and the whole APCM becomes as shown in **Algorithm 3**. Here, l is a N -dimensional vector, whose i th component contains the index of the cluster which is most compatible with the data vector \mathbf{x}_i (i.e. the cluster C_j for which $u_{ij} = \max_{r=1, \dots, m} u_{ir}$), n_j denotes the number of the data points \mathbf{x}_i , that are most compatible with the j th cluster and μ_j is the mean vector of these data points. In words, the "Possible cluster elimination" part of APCM examines if the index j of a cluster C_j appears in the vector l . If this is the case (i.e. if there exists at least one vector \mathbf{x}_i that is most compatible with C_j), C_j is not eliminated. Otherwise, C_j is eliminated. Also, η_j is estimated as the variance of the most compatible with C_j points, around their mean μ_j .

Let us first comment on the proposed updating mechanism of η_j 's. Note that, although their updating formula resembles to that used in the classical PCM, as well as many of its variants, it differs from them in two distinctive points. First η_j 's are updated taking into account only the data vectors that are most compatible to C_j . Second, the distances involved in the formula are between a data vector and the mean vector $\mu_j(t)$, not the representative $\theta_j(t)$. This allows more accurate estimates for η_j 's. It is also noted that, in the case where there are two or more clusters, that are equally compatible with a specific \mathbf{x}_i , then \mathbf{x}_i will contribute to the determination of the η parameter of only one (arbitrarily chosen) of them.

Let us now comment on the immunity of the proposed algorithm to overestimates on the number of clusters. In such a scenario and taking into account that a possibilistic type algorithm moves the representatives in dense regions, the prob-

¹ *.*\$4(*33:"/054 4*/y # \$- .#*.'02/"\$ #5\$2350 \$"y %12.2' \$: # 43:\$432\$15*2\$3.Q #0\$21 4*/.6\$:#/04\$2\$1(\$002/7*- 4.\$: ,/6:"/34;'/2*4(#\$3"2*!\$#:

Algorithm 3 The APCM algorithm

```

Initialize  $\mathbf{M} \leftarrow \mathbf{M}^0$  i f p ^ P @ { f m u p f l }
Initialize  $\mathbf{Q} \leftarrow \mathbf{Q}^0$  i @ { m u p f l }
t = 0
Repeat:
  Update  $U$  (as in Algorithm 1)
  t = t + 1
  Update  $\Theta$  (as in Algorithm 1)
  Possible cluster identification
  for i = 1 to N
    Determine  $\mathbf{1}; \mathbf{B} = \mathbf{B}^* \mathbf{M} \mathbf{1}; \mathbf{C} = \mathbf{C}^*$ 
    Set:  $l a b e l r(i) =$ 
  end for
  for j = 1 to m
    If  $j \notin a$  then
      Remove  $\mathbf{C}^*$ 
       $m = m - 1$ 
    end if
  end for
  Adaptation
  for j = 1 to m
     $\mathbf{Q}(t) = \frac{1}{n_j(t)} \sum_{\mathbf{x}_i: u_{ij}(t) = \max_{r=1, \dots, m} u_{ir}(t)} |\mathbf{x}_i - \mathbf{Q}_j(t)|$ 
  end for
until: a termination criterion is met

```

ability to select as representative at least one point in each dense region is increased, due to the way the representatives are initially selected (via the Max-Min Algorithm). Then, due to the possibilistic nature of APCM, it is guaranteed, in principle, that the number of the representatives which move to a specific dense region will be reduced to a single one.

In order to get some further insight to the way the algorithm works, assume that two cluster representatives \mathbf{B}_i and \mathbf{B}_j most coincide but let say $\mathbf{B}_i \neq \mathbf{B}_j$. Then, for a given data point \mathbf{x}_i , it is $\frac{\|\mathbf{x}_i - \mathbf{B}_i\|}{\eta_i} < 1 < \frac{\|\mathbf{x}_i - \mathbf{B}_j\|}{\eta_j}$, which implies that $\mathbf{1} \in \mathbf{M}_i$ considering the updating equation of $\mathbf{1}; \mathbf{C}$ in the possibilistic algorithm. Loosely speaking, between two cluster representatives with different parameters, the one with the greater η has a stronger influence around it. Thus, if \mathbf{B}_i and \mathbf{B}_j most coincide, the influence of the one with the smaller η will be vanished by the influence of the one with the greater η in the sense that $\mathbf{1}; \mathbf{B}_i$, for all data points \mathbf{x}_i .

4 EXPERIMENTAL RESULTS

In this section, we test the proposed method through a synthetic data set and two real data sets. Moreover, we compare the results with those obtained from k-means, FCM and PCM. Note that, because the PCM algorithm leads, in some cases, to coincident clusters, the clustering result is extracted taking

Table 1 The results of the synthetic data set

	$m_{initial}$	m_{final}	Rand Measure
k-means	3	3	97.2%
k-means	5	5	82.57%
FCM	3	3	96.8%
FCM	5	5	82.20%
PCM	3	2	83.98%
PCM	5	2	83.98%
PCM	10	2	83.98%
APCM	5	3	96.9%
APCM	10	3	96.9%

into account only the truly “different” clusters. In addition, in order to make a fair comparison between all clustering algorithms, the representatives \mathbf{M}^0 initialized based on the Max-Min scheme, in all of them.

4.1 A synthetic data set experiment

Let us consider a two-dimensional data set consisting of $N = 1000$ points that form three clusters. Each cluster is modelled by a normal distribution. Their means are $\mathbf{Q}_1 = [1.0735, 0.75]^T$, $\mathbf{Q}_2 = [7.3807, 3.34]^T$ and $\mathbf{Q}_3 = [5.34, 3.34]^T$, respectively, while their (common) covariance matrix is $0.8 \cdot I_2$, where I_2 is the 2×2 identity matrix. A number of 400 points have been generated by each one of the first and the third distribution, while 200 points have been generated by the second one. Note that the second and the third clusters are close enough to each other, while they are far away from the first one. Each data point is assigned to a cluster, utilizing the U matrix, as follows: \mathbf{x}_i is assigned to cluster k if $u_{ik} > \max_{r \neq k} u_{ir}$. Figs. 1 (a), (b) show the clustering results obtained using the k-means algorithm with $m = 3$ and $m = 5$, respectively. Similarly, in Figs. 1 (c), (d) we present the corresponding results for FCM. Fig. 1 (e) depicts the performance of PCM for $m = 5$, while in addition, it shows the circled regions, centered at each \mathbf{Q}_j and having radius equal to η_j , in which \mathbf{M}^0 has increased influence. Finally, Fig. 1 (f) shows the results of APCM with $m = 5$ and $k = 0, 1$.

In order to compare a clustering with the true data label information, we use the Rand Measure described in [1]. Table 1 shows the results of the previously mentioned algorithms for the synthetic data set, where $m_{initial}$ and m_{final} denote the initial and the final number of clusters, respectively. It presents the Rand Measure of each algorithm, which results from the comparison of the $l a b e l r$, and the true cluster identity of the vectors.

As it is deduced from Fig. 1 and Table 1, when k-means and FCM are initialized by the (unknown in practice) true number of clusters ($m = 3$) the provided results are very satisfactory. However, any deviation from this value causes a significant degradation to the obtained clustering quality. On the other hand, the classical PCM fails to unravel the under-

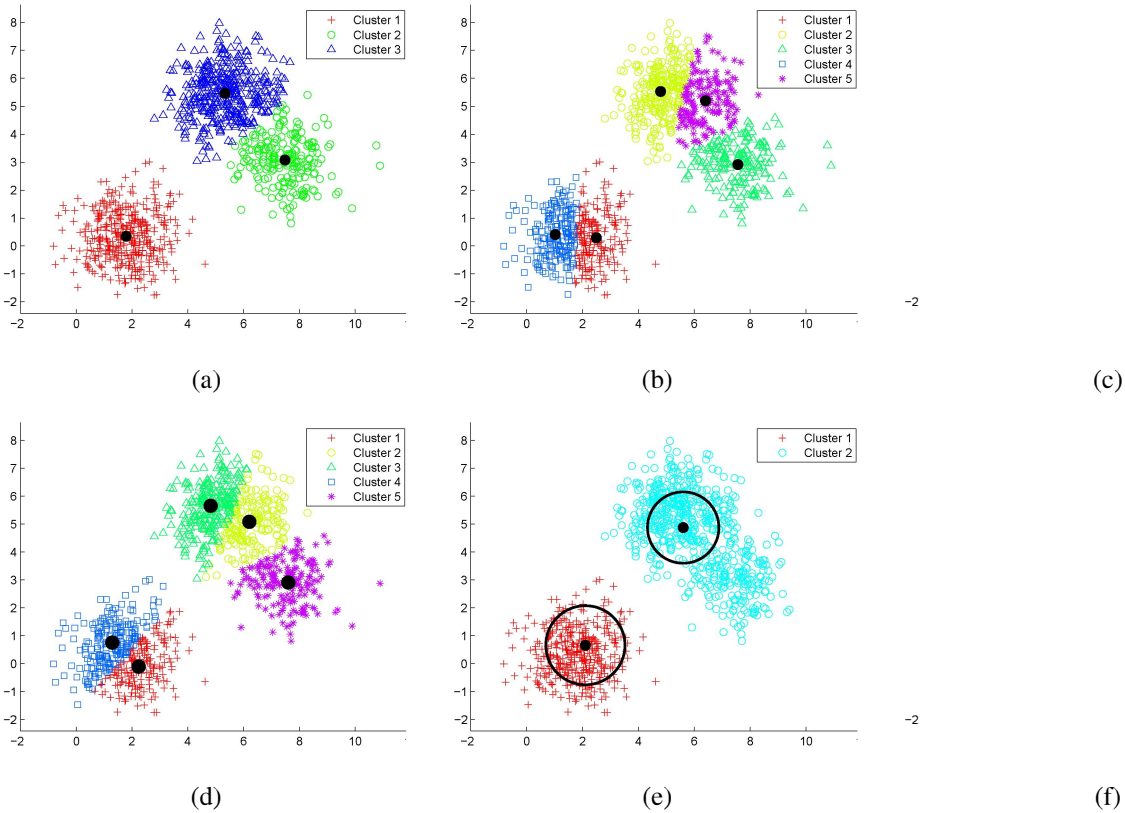


Fig 1 Clustering results of k-means with, (a) $m = 3$ and (b) $m = 5$ FCM with, (c) $m = 3$ and (d) $m = 5$ (e) PCM with $m = 5$ and (f) APCM with $m = 5$. Note that in PCM the clustering result is extracted taking into account only the truly “different” clusters. Bolded dots represent the final clusters’ representatives.

ying clustering structure, due to the fact that two clusters are close enough to each other and the algorithm does not have the ability to adapt η_j ’s in order to cope with this situation. Finally, the proposed APCM constantly produces very accurate results for various initial values of m .

4.2 Real data experiments

In this sub-section, the Iris data set and the Wine data set from UCI Library database [11] are considered. In order to efficiently utilize all the information of the l features of these real data sets, we shall normalize them as

$$\hat{x}_{ip} = \frac{x_{ip} - \frac{1}{N} \sum_{p=1}^N x_{rp}}{\sqrt{\frac{1}{N-1} \sum_{p=1}^N (x_{rp} - \frac{1}{N} \sum_{p=1}^N x_{rp})^2}}$$

where $i = 1, \dots, N$, $p = 1, \dots, l$.

Iris data set This set consists of $N = 150$ 4-dimensional data points that form three classes, each one having 50 points. In this data set, two classes are overlapped thus one can argue whether the number of clusters m is 2 or 3.

As it is shown in Table 2, APCM ends up with $m_{final} = 2$ clusters independently of the initial number of clusters $m_{initial}$. Since algorithms following the possibilistic philosophy detect dense regions, the APCM ends with $m_{final} = 2$ clusters, which is an acceptable result, due to the nature of the data set. On the other hand, k-means and FCM perform better provided that they have been supplied with the true number of the underlying classes. However, this situation changes as we deviate from this value. In addition, it is interesting to note that for $m_{initial} = 2$ both k-means and FCM provide the same clustering with that produced by APCM. This is an indication that APCM provides constantly the best possible two-cluster clustering for all initial values for m and $\beta = 0.1$. Finally, the original PCM needs a significantly greater than 3 initial value for m in order to be competitive with the other schemes.

Wine data set: This set consists of $N = 178$, 3-dimensional data points that stem from three classes, the first with 59 points, the second with 71 and the third one with 48 points.

In the Wine data set, an exhaustive feature selection procedure preceded the clustering stage, due to the high dimensionality of this data set in relation to the small number of data

Table 2 The results of the real Iris data set

	t	t	η	Final β	Final Measure
k-means	2	2	2	2	76.61%
k-means	3	3	3	3	75.38%
k-means	5	5	5	5	75.38%
k-means	8	8	8	8	75.38%
FCM	2	2	2	2	77.43%
FCM	3	3	3	3	77.43%
FCM	5	5	5	5	77.43%
FCM	8	8	8	8	77.43%
PCM	2	1	1	1	32.89%
PCM	3	1	1	1	32.89%
PCM	5	1	1	1	32.89%
PCM	8	2	2	2	74.98%
APCM	4	2	2	2	77.43%
APCM	5	2	2	2	77.43%
APCM	8	2	2	2	77.43%
APCM	10	2	2	2	77.43%

Table 3 The results of the real Wine data set

	t	t	η	Final β	Final Measure
k-means	3	3	3	3	82.30%
k-means	5	5	5	5	89.19%
k-means	8	8	8	8	82.30%
FCM	3	3	3	3	93.7%
FCM	5	5	5	5	83.75%
FCM	8	8	8	8	78.22%
PCM	3	1	1	1	33.80%
PCM	5	1	1	1	33.80%
PCM	8	1	1	1	33.80%
APCM	4	3	3	3	93.7%
APCM	5	3	3	3	93.7%
APCM	8	3	3	3	93.7%

points. In order to make a fair comparison between the clustering algorithms, as presented in Table. 3, we considered the optimal combination of features for each algorithm. It turned out that the best result for each algorithm was obtained using different combinations of eight features. As before, Table 3 shows that APCM ($\beta = 0.1$) results in ~ 3 clusters independently of the initialization of the number of clusters t while the other clustering algorithms require the knowledge of the actual number of clusters, in order to extract satisfactory results.

5 CONCLUSIONS

In this paper a novel possibilistic clustering algorithm (APCM) has been proposed. The algorithm encompasses a proper initialization and a new updating mechanism for the η parameters and is immune to overestimates on the actual number of existing clusters. These features make the algorithm very

flexible in tracking the clustering environment under study. The performance of the proposed algorithm against k-means, FCM and the original PCM has been assessed using both synthetic and real data sets. In all these experiments, it is shown that APCM has a steadily good performance irrespective of the initial number of clusters, which is not the case with k-means and FCM. In addition, APCM constantly provides better results compared to the original PCM.

6 REFERENCES

- [1] S. Theodoridis and K. Koutroumbas, *Pattern Recognition in the Real World*, Academic Press, 2009.
- [2] J. C. Bezdek, *Pattern Recognition with Fuzzy Objective Functions*, John Wiley & Sons, 1981.
- [3] R. Krishnapuram and J. M. Keller, "A possibilistic approach to clustering," *IEEE Transactions on Fuzzy Systems*, vol. 1, pp. 98–110, 1993.
- [4] M. Barni, V. Cappellini, and A. Mecocci, "A possibilistic approach to clustering," *IEEE Transactions on Fuzzy Systems*, vol. IV, pp. 393–396, 1996.
- [5] M.-S. Yang and K.-L. Wu, "Unsupervised possibilistic clustering," *Pattern Recognition*, 39, pp. 5–21, 2006.
- [6] K. Treerattanapitak and C. Jaruskulchai, "Possibilistic exponential fuzzy clustering," *Journal of Computer Science*, vol. 28, pp. 311–321, 2013.
- [7] J.-S. Zhang and Y.-W. Leung, "Improved possibilistic c-means clustering algorithms," *IEEE Transactions on Fuzzy Systems*, vol. 12, pp. 209–217, 2004.
- [8] R. Krishnapuram and J. M. Keller, "The possibilistic c-means algorithm: insights and recommendations," *IEEE Transactions on Fuzzy Systems*, vol. IV, pp. 385–393, 1996.
- [9] B. Mirkin, *Clustering Data Analysis*, John Wiley & Sons, 2005.
- [10] O. Egecioglu and B. Kalantari, "Approximating the diameter of a set of points in the euclidean space," *Information Science*, vol. 32, pp. 205–211, 1989.
- [11] UCI Machine Learning Database, <http://archive.uci.edu/ml/dataset.html>