

# Fuzzy Cluster Analysis with Cluster Repulsion

Heiko Timm, Christian Borgelt, and Rudolf Kruse  
Dept. of Knowledge Processing and Language Engineering  
Otto-von-Guericke-University of Magdeburg  
Universitätsplatz 2, D-39106 Magdeburg, Germany  
{timm,borgelt,kruse}@iws.cs.uni-magdeburg.de

## Abstract

We explore an approach to possibilistic fuzzy  $c$ -means clustering that avoids a severe drawback of the conventional approach, namely that the objective function is truly minimized only if all cluster centers are identical. Our approach is based on the idea that this undesired property can be avoided if we introduce a mutual repulsion of the clusters, so that they are forced away from each other. In our experiments we found that in this way we can combine the partitioning property of the probabilistic fuzzy  $c$ -means algorithm with the advantages of a possibilistic approach w.r.t. the interpretation of the membership degrees.

## 1 Introduction

Cluster analysis is a technique for classifying data, i.e., to divide a given dataset into a set of classes or *clusters*. The goal is to divide the dataset in such a way that two cases from the same cluster are as similar as possible and two cases from different clusters are as dissimilar as possible. Thus one tries to model the human ability to group similar objects or cases into classes and categories. In classical cluster analysis each datum must be assigned to exactly one cluster. Fuzzy cluster analysis relaxes this requirement by allowing gradual memberships, thus offering the opportunity to deal with data that belong to more than one cluster at the same time.

Most fuzzy clustering algorithms are objective function based: They determine an optimal classification by minimizing an objective function. In objective function based clustering usually each cluster is represented by a *cluster prototype*. This prototype consists of a *cluster center* (whose name already indicates its meaning) and maybe some additional information about the size and the shape of the cluster. The cluster center is an instantiation of the attributes used to describe the domain, just as the data points in the dataset to divide. However, the cluster center is computed by the clustering algorithm and may or may not appear in the dataset. The size and shape parameters determine the extension of the cluster in different directions of the underlying domain.

The degrees of membership to which a given data point belongs to the different clusters are computed from the distances of the data point to the cluster centers w.r.t. the size and the shape of the cluster as stated by the additional prototype information. The closer a data point lies to the center of a cluster (w.r.t. size and shape), the higher is its degree of membership to this cluster. Hence the problem to divide a dataset  $X = \{\vec{x}_1, \dots, \vec{x}_n\} \subseteq \mathbb{R}^p$  into  $c$  clusters can be stated as the task to minimize the distances of the data points to the cluster centers, since, of course, we want to maximize the degrees of membership.

Several fuzzy clustering algorithms can be distinguished depending on the additional size and shape information contained in the cluster prototypes, the way in which the distances are determined, and the restrictions that are placed on the membership degrees. Here we focus on the fuzzy  $c$ -means algorithm [1], which uses only cluster centers and a Euclidean distance function. We distinguish, however, between probabilistic and possibilistic clustering, which use different sets of constraints for the membership degrees.

## Probabilistic Fuzzy Clustering

In probabilistic fuzzy clustering the task is to minimize the objective function

$$J(\mathbf{X}, \mathbf{U}, \mathbf{B}) = \sum_{i=1}^c \sum_{j=1}^n u_{ij}^m d^2(\vec{\beta}_i, \vec{x}_j) \quad (1)$$

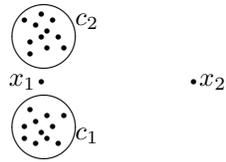


Figure 1: A situation in which the probabilistic assignment of membership degrees is counterintuitive for datum  $x_2$ .

subject to

$$\sum_{j=1}^n u_{ij} > 0, \quad \text{for all } i \in \{1, \dots, c\}, \quad \text{and} \quad (2)$$

$$\sum_{i=1}^c u_{ij} = 1, \quad \text{for all } j \in \{1, \dots, n\}, \quad (3)$$

where  $u_{ij} \in [0, 1]$  is the membership degree of datum  $\vec{x}_j$  to cluster  $c_i$ ,  $\vec{\beta}_i$  is the prototype of cluster  $c_i$ , and  $d(\vec{\beta}_i, \vec{x}_j)$  is the distance between datum  $\vec{x}_j$  and prototype  $\vec{\beta}_i$ .  $\mathbf{B}$  is the set of all  $c$  cluster prototypes  $\vec{\beta}_1, \dots, \vec{\beta}_c$ . The  $c \times n$  matrix  $\mathbf{U} = [u_{ij}]$  is called the fuzzy partition matrix and the parameter  $m$  is called the fuzzifier. This parameter determines the “fuzziness” of the classification. With higher values for  $m$  the boundaries between the clusters become softer, with lower values they get harder. Usually  $m = 2$  is chosen.

Constraint (2) guarantees that no cluster is empty and constraint (3) ensures that the sum of the membership degrees for each datum equals 1. Because of the second constraint, this approach is called *probabilistic clustering*, since with it the membership degrees for a given datum formally resemble the probabilities of its being a member of the corresponding cluster.

Unfortunately, the objective function  $J$  cannot be minimized directly. Therefore an iterative algorithm is used, which alternately optimizes the cluster prototypes and the membership degrees. That is, first the cluster prototypes are optimized for fixed membership degrees, then the membership degrees are optimized for fixed prototypes. The main advantage of this scheme is that in each of the two steps the optimum can be computed directly. By iterating the two steps the joint optimum is approached. The update formulae are derived by simply setting the derivative of the objective function (extended by Lagrange multipliers to incorporate the constraints) w.r.t. the parameter to optimize equal to zero. For the membership degrees we thus obtain the following formula

$$u_{ij} = \begin{cases} \frac{1}{\sum_{k=1}^c \left( \frac{d^2(x_j, \beta_i)}{d^2(x_j, \beta_k)} \right)^{\frac{1}{m-1}}}, & \text{if } I_j = \emptyset, \\ 0, & \text{if } I_j \neq \emptyset \text{ and } i \notin I_j, \\ x, x \in [0, 1] \text{ such that } \sum_{i \in I_j} u_{ij} = 1, & \text{if } I_j \neq \emptyset \text{ and } i \in I_j. \end{cases} \quad (4)$$

Equation (4) shows that the membership degree of a datum to a cluster depends not only on the distance between the datum and that cluster, but also on the distances between the datum and other clusters. The partitioning property of a probabilistic clustering algorithm, which “distributes” the weight of a datum on the different clusters, is due to this equation.

Although often desirable, the “relative” character of the membership degrees in a probabilistic clustering approach can lead to counterintuitive results. Consider, for example, the simple case of two clusters shown in figure 1. Datum  $\vec{x}_1$  has the same distance to both clusters and thus it is assigned a degree of membership of about 0.5. This is plausible. However, the same degrees of membership are assigned to datum  $\vec{x}_2$ . Since this datum is far away from both clusters, it would be more intuitive if it had a low degree of membership to both of them.

## Possibilistic Fuzzy Clustering

In possibilistic fuzzy clustering one tries to achieve a more intuitive assignment of degrees of membership by dropping constraint (3), which is responsible for the undesirable effect discussed above. However, this leads to the mathematical problem that the objective function is now minimized by assigning  $u_{ij} = 0$  for all  $i \in \{1, \dots, c\}$  and  $j \in \{1, \dots, n\}$ . In order to avoid this trivial solution, a penalty term is introduced, which forces the

membership degrees away from zero. That is, the objective function  $J$  is modified to

$$J(\mathbf{X}, \mathbf{U}, \mathbf{B}) = \sum_{i=1}^c \sum_{j=1}^n u_{ij}^m d^2(\vec{\beta}_i, \vec{x}_j) + \sum_{i=1}^c \eta_i \sum_{j=1}^n (1 - u_{ij})^m, \quad (5)$$

where  $\eta_i > 0$ . The first term leads to a minimization of the weighted distances while the second term suppresses the trivial solution. This approach is called *possibilistic clustering*, because the membership degrees for one datum resemble the possibility (in the sense of possibility theory [6]) of its being a member of the corresponding cluster [10, 5]. The formula for updating the membership degrees that is derived from this objective function is [10]

$$u_{ij} = \frac{1}{1 + \left( \frac{d^2(\vec{x}_j, \vec{\beta}_i)}{\eta_i} \right)^{\frac{1}{m-1}}}. \quad (6)$$

From this equation it becomes obvious that  $\eta_i$  is a parameter that determines the distance at which the membership degree equals 0.5.  $\eta_i$  is chosen for each cluster separately and can be determined, for example, by computing the fuzzy intra cluster distance [10]

$$\eta_i = \frac{K}{N_i} \sum_{j=1}^n u_{ij}^m d^2(\vec{x}_j, \vec{\beta}_i), \quad (7)$$

where  $N_i = \sum_{j=1}^n u_{ij}^m$ . Usually  $K = 1$  is chosen.

At first sight this approach looks very promising. However, if we take a closer look, we discover that the objective function  $J$  defined above is, in general, truly minimized only if all cluster centers are identical. The reason is that the formula (6) for the membership degree of a datum to a cluster depends only on the distance of the datum to that cluster, but not on its distance to other clusters. Hence, if there is a single optimal point for a cluster center (as it will usually be the case, since multiple optimal points would require a high symmetry in the data), all cluster centers will be moved there. More formally, consider two cluster centers  $\vec{\beta}_1$  and  $\vec{\beta}_2$ , which are not identical, and let

$$z_i = \sum_{j=1}^n u_{ij}^m d^2(\vec{\beta}_i, \vec{x}_j) + \eta_i \sum_{j=1}^n (1 - u_{ij})^m, \quad i = 1, 2,$$

i.e., let  $z_i$  be the amount that cluster  $\beta_i$  contributes to the value of the objective function. Except in very rare cases of high data symmetry, it will then either be  $z_1 > z_2$  or  $z_2 > z_1$ . That is, we can improve the value of the objective function by setting both cluster centers to the same value, namely the one which yields the smaller  $z$ -value, because the two  $z$ -values do not interact.

Note that this behavior is specific to the possibilistic approach. In the probabilistic approach the cluster centers are driven apart, because a cluster, in a way, “seizes” part of the weight of a datum and thus leaves less that may attract other cluster centers. Hence sharing a datum between clusters is disadvantageous. In the possibilistic approach there is nothing to complement this effect.

Nevertheless, possibilistic fuzzy clustering usually leads to acceptable results, although it suffers from stability problems if it is not initialized with the corresponding probabilistic algorithm. We assume that other results than all cluster centers being identical are achieved only, because the algorithm gets stuck in a local minimum of the objective function. This, of course, is not a desirable situation. Hence we tried to improve the algorithm by modifying the objective function in such a way that the problematic property examined above is removed.

## 2 A New Approach Based on Cluster Repulsion

The idea of our approach is to combine an attraction of data to clusters with a repulsion between different clusters. In contrast to a probabilistic clustering algorithm this is not done implicitly using restriction (3), but explicitly by adding a cluster repulsion term to the objective function.

To arrive at a suitable objective function, we started from the following set of requirements:

- The distance between clusters and the data points assigned to them should be minimized.
- The distance between clusters should to be maximized.

- There should be no empty clusters, i.e., for each cluster there must be datum with non-vanishing membership degree.
- Membership degrees should be close to one and, of course, the trivial solution of all membership degrees being zero should be suppressed.

These requirements are very close to standard possibilistic cluster analysis. The attraction between data and clusters is modeled (as described above) by a term  $\sum_{i=1}^c \sum_{j=1}^n u_{ij}^m d^2(\vec{\beta}_i, \vec{x}_j)$ . A term  $\sum_{i=1}^c \eta_i \sum_{j=1}^n (1 - u_{ij})^m$  is used to avoid the trivial solution. The objective that to each cluster data have to be assigned is leads to the constraint (2). The repulsion between clusters can be described in analogy to the attraction between data and clusters. That is, we are using a term that is minimized if the sum of the distances between clusters are maximized.

This could be achieved by simply subtracting the sum of squared distances between clusters from the objective function. However, this straightforward approach does not work. The problem is that the repulsion then increases with the distance of the clusters and thus driving them ever farther apart improves the value of the objective function. In the end, all data points would be assigned to one cluster and all other clusters would have been moved to infinity.

To avoid this undesired “explosion” of the cluster set, a repulsion term must be used that gets smaller the farther the clusters are apart. Then the attraction of the data points can compensate the repulsion if only the clusters are sufficiently spread out. This consideration lead us to the term  $\gamma \sum_{i=1}^c \sum_{k=1, k \neq i}^c \frac{1}{d^2(\vec{\beta}_i, \vec{\beta}_k)}$  where  $\gamma$  is a weighting factor. This term is only relevant if the clusters are close together. With growing distance it becomes smaller, i.e., the repulsion is gradually decreased until it is compensated by the attraction of the data.

The classification problem is then described as the task to minimize

$$J(\mathbf{X}, \mathbf{U}, \mathbf{B}) = \sum_{i=1}^c \sum_{j=1}^n u_{ij}^m d^2(\vec{\beta}_i, \vec{x}_j) + \sum_{i=1}^c \eta_i \sum_{j=1}^n (1 - u_{ij})^m + \gamma \sum_{i=1}^c \sum_{k=1, k \neq i}^c \frac{1}{d^2(\vec{\beta}_i, \vec{\beta}_k)} \quad (8)$$

w.r.t. the constraint  $\sum_{j=1}^n u_{ij} > 0$  for all  $i \in \{1, \dots, c\}$ .  $\gamma$  is used to weight the objective that the distance to the clusters should be minimized against the objective that the distance between clusters should be maximized. Using  $\frac{1}{d^2(\vec{\beta}_i, \vec{\beta}_k)}$  means that only clusters with a small distance are relevant for minimizing the objective function, while clusters with a large distance are only slightly repelling each other.

Minimization of (8) w.r.t. the membership degrees leads to (6). That is, the membership degrees have the same meaning as in possibilistic cluster analysis. For the variant of the fuzzy  $c$ -means algorithm (only cluster centers  $\vec{c}_i$ , Euclidean distance, and therefore spherical clusters) a minimization of (8) with respect to the cluster prototypes leads to

$$\sum_{j=1}^n u_{ij} (\vec{x}_j - \vec{c}_i) - \sum_{k=1, k \neq i}^c (\vec{c}_k - \vec{c}_i) \frac{1}{\|\vec{c}_k - \vec{c}_i\|^2} = 0. \quad (9)$$

For reasons of simplicity, we solved (9) by iteratively computing

$$\vec{c}_i = \frac{\sum_{j=1}^n u_{ij} \vec{x}_j - \gamma \sum_{k=1, k \neq i}^c \vec{c}_k \frac{1}{\|\vec{c}_k - \vec{c}_i\|^2}}{\sum_{j=1}^n u_{ij} - \gamma \sum_{k=1, k \neq i}^c \frac{1}{\|\vec{c}_k - \vec{c}_i\|^2}} \quad (10)$$

For  $\vec{c}_i$  on the right hand side we used old values of the previous iteration. The computation was iterated until  $|\vec{c}_i^{(\text{new})} - \vec{c}_i^{(\text{old})}| < \bar{\epsilon}$ .

(10) shows the effect of the repulsion between clusters. A cluster is attracted by the data assigned to it and repelled by the other clusters.

An alternative approach to model the repulsion between clusters is to use the term  $\gamma \sum_{i=1}^c \sum_{k=1, k \neq i}^c e^{-d^2(\vec{\beta}_i, \vec{\beta}_k)}$  instead of the fraction used above. The difference between both terms is how the repulsion between clusters decreases with a growing distance.

The classification problem is then described as the task to minimize

$$J(\mathbf{X}, \mathbf{U}, \mathbf{B}) = \sum_{i=1}^c \sum_{j=1}^n u_{ij}^m d^2(\vec{\beta}_i, \vec{x}_j) + \sum_{i=1}^c \eta_i \sum_{j=1}^n (1 - u_{ij})^m + \gamma \sum_{i=1}^c \sum_{k=1, k \neq i}^c e^{-d^2(\vec{\beta}_i, \vec{\beta}_k)} \quad (11)$$

w.r.t. the constraint  $\sum_{j=1}^n u_{ij} > 0$  for all  $i \in \{1, \dots, c\}$ .

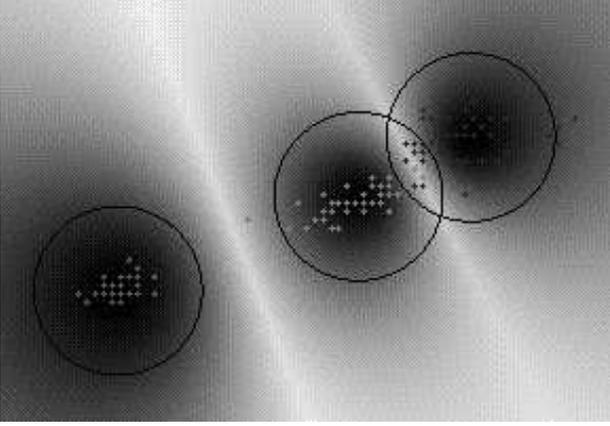


Figure 2: Iris dataset classified with probabilistic fuzzy  $c$ -means algorithm. Attributes petal length and petal width.

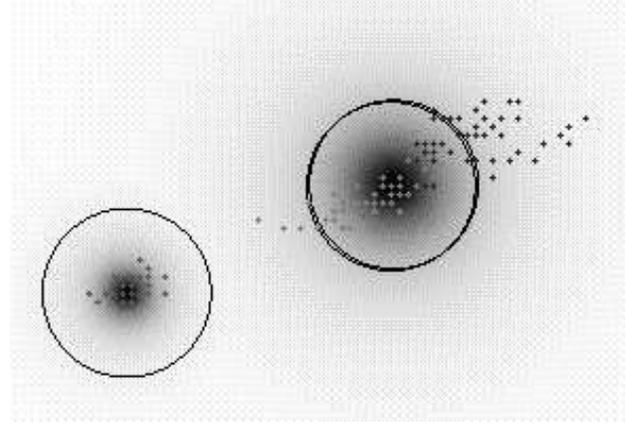


Figure 3: Iris dataset classified with possibilistic fuzzy  $c$ -means algorithm. Attributes petal length and petal width.

Minimizing (11) w.r.t.  $\vec{\beta}_i$  leads for the fuzzy  $c$ -means algorithm, that is, if the clusters are described by their centers  $\vec{c}_i$  only, to

$$\sum_{j=1}^n u_{ij}(\vec{x}_j - \vec{c}_i) - \sum_{k=1, k \neq i}^c (\vec{c}_k - \vec{c}_i)e^{-\|\vec{c}_k - \vec{c}_i\|} = 0. \quad (12)$$

As (9) we solved (12) by an iterative approach.

In the approaches presented in this section the attraction between clusters and data assigned to them and the repulsion between clusters is modeled separately. In contrast to a probabilistic clustering algorithm the membership degree can be interpreted as a measure of similarity to a cluster. The repulsion between clusters avoids the problems of possibilistic cluster analysis as described above.  $\gamma$  is used to weight the two opposite objectives, i.e., that the distance between clusters and data assigned to them should be minimized and that the distance between clusters should be maximized.

### 3 Test Examples

We used the well-known iris data set [7] for testing our algorithm. We used only the attributes petal length and petal width, since these carry the most information about the distribution of the iris flowers. Fig. 2 shows the classification obtained with the probabilistic fuzzy  $c$ -means algorithm. This result clearly demonstrates the partitioning property of the probabilistic algorithm. The data set is divided into three clusters. Fig. 3 shows the classification obtained with the possibilistic fuzzy  $c$ -means algorithm. Only two clusters are detected because the possibilistic algorithm is not forced to partition the data. As shown in section 1 the two clusters on the right are almost identical. The cluster on the left is detected, because it is well separated and thus forms a local minimum of the objective function.

Fig. 4, 5, 6, and 7 show the results of minimizing the objective function 8 and fig. 8, 9, 10, and 11 the results of minimizing the objective function 11 for different values of  $\gamma$ . The classification is computed using possibilistic membership degrees as described in section 2. However, in contrast to standard possibilistic cluster analysis, three clusters are detected. Using cluster repulsion leads to a classification similar to the result of probabilistic clustering. We computed the classification with several values for  $\gamma$ . The method seems to be very robust with respect to the choice of the weighting factor  $\gamma$ .

### 4 Conclusion and Future Work

In this paper we presented an approach for possibilistic fuzzy cluster analysis that is based on data attracting cluster centers as well as cluster centers repelling each other. This approach combines the more intuitive membership degrees of possibilistic fuzzy cluster analysis (since they can be interpreted as similarities) with the partitioning property of probabilistic cluster analysis. By this we combine the advantages of both approaches.

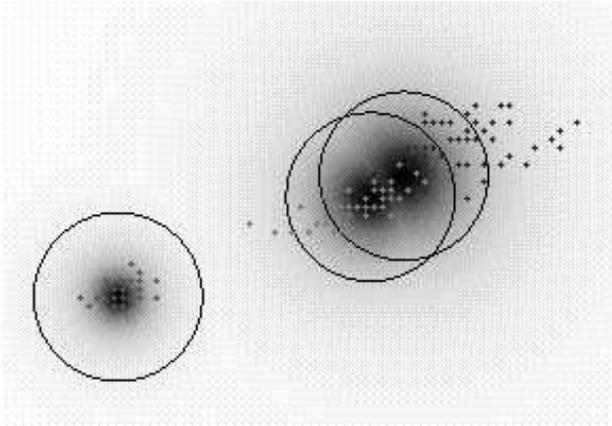


Figure 4: Iris dataset classified with approach based on objective function (8).  $\gamma = 0.1$ . Attributes petal length and petal width.

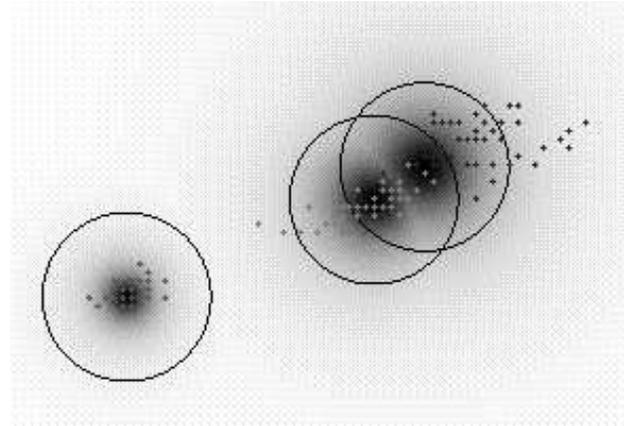


Figure 5: Iris dataset classified with approach based on objective function (8).  $\gamma = 0.5$ . Attributes petal length and petal width.

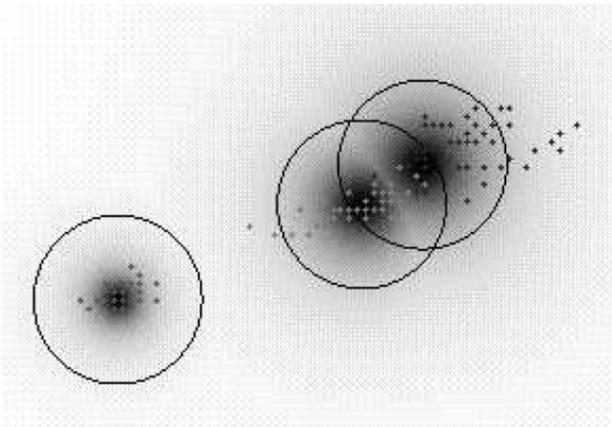


Figure 6: Iris dataset classified with approach based on objective function (8).  $\gamma = 1$ . Attributes petal length and petal width.

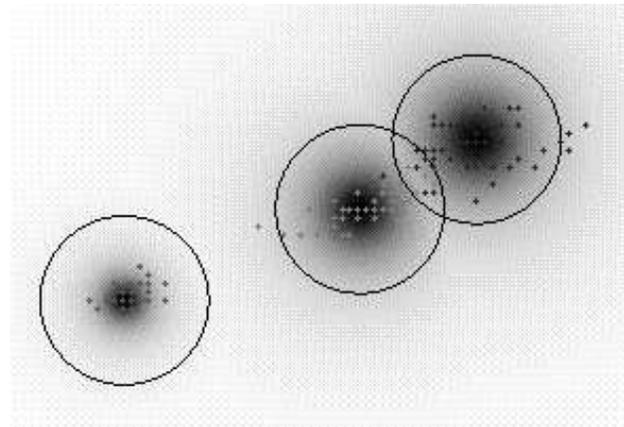


Figure 7: Iris dataset classified with approach based on objective function (8).  $\gamma = 10$ . Attributes petal length and petal width.

In the future we plan to extend the approach presented in this paper to other fuzzy clustering algorithms as, for instance, the Gustafson-Kessel algorithm. Furthermore we plan to study how to extend it to deal with classified data. In [11] this was done using a repulsion between data and clusters belonging to different classes. However, this can also be done by a possibilistic clustering algorithm as described in this paper with weights  $\gamma_{\text{equal class}}$  and  $\gamma_{\text{different classes}}$ . Another idea would be to use a probabilistic fuzzy clustering algorithm with a repulsion between clusters belonging to different classes as described in this paper.

## References

- [1] Bezdek, J.C.: Pattern Recognition with Fuzzy Objective Function Algorithms. Plenum Press, New York, NY, USA 1981.
- [2] Bezdek, J.C., Keller, J., Krishnapuram R., and Pal, N.R.: Fuzzy Models and Algorithms for Pattern Recognition and Image Processing. Kluwer, Boston, London, 1999.
- [3] Bezdek, J.C. and Pal S.K.: Fuzzy Models for Pattern Recognition — Methods that Search for Structures in Data. IEEE Press, Piscataway, NJ, USA 1992.
- [4] Borgelt, C. bcview: A program to visualize the numeric part of full or a naive Bayes classifier. <http://fuzzy.uni-magdeburg.de/~borgelt/software.html>.

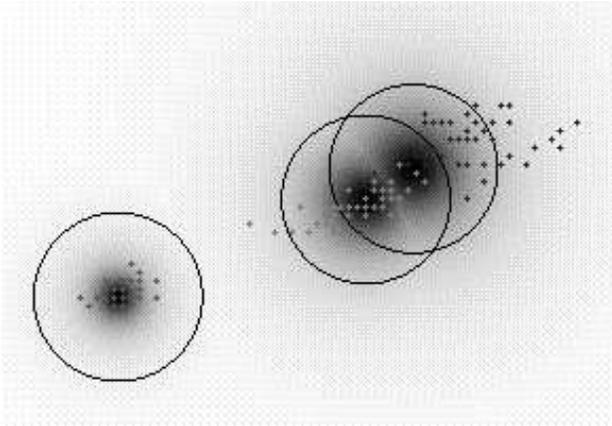


Figure 8: Iris dataset classified with approach based on objective function (11).  $\gamma = 3$ . Attributes petal length and petal width.

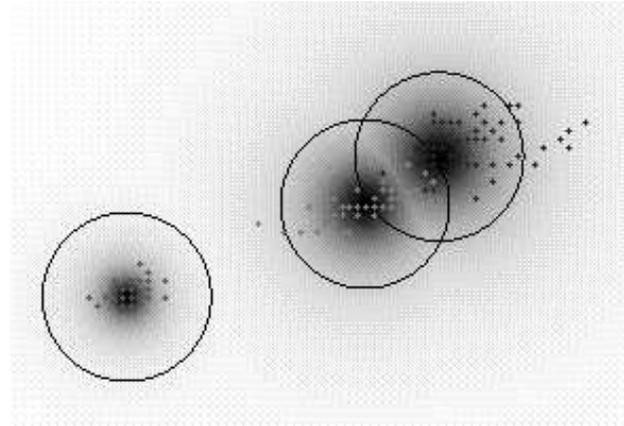


Figure 9: Iris dataset classified with approach based on objective function (11).  $\gamma = 5$ . Attributes petal length and petal width.

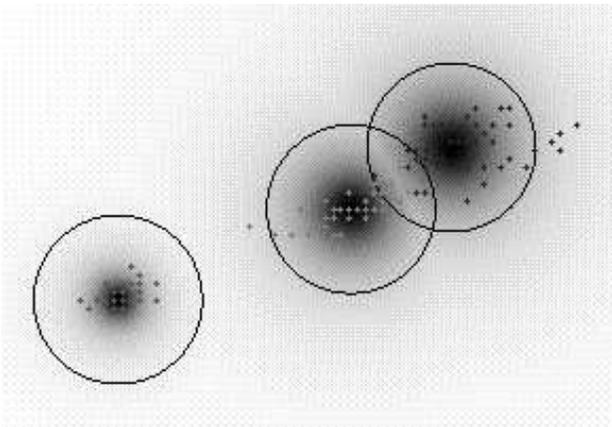


Figure 10: Iris dataset classified with approach based on objective function (11).  $\gamma = 10$ . Attributes petal length and petal width.

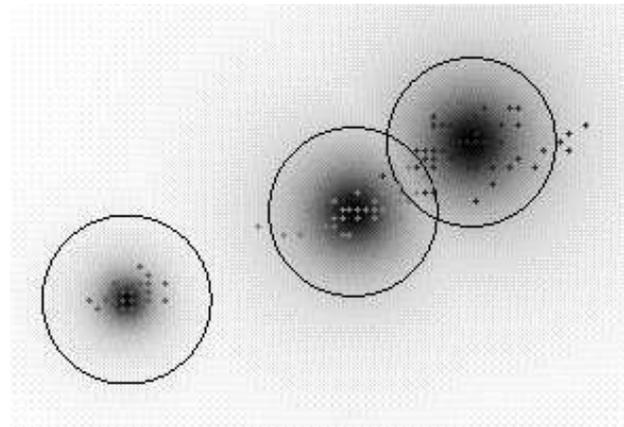


Figure 11: Iris dataset classified with approach based on objective function (11).  $\gamma = 20$ . Attributes petal length and petal width.

- [5] Davé, R.N. und Krishnapuram, R.: Robust Clustering Methods: A Unified View, *IEEE Transactions on Fuzzy Systems*, pp. 270-293, (5) 1997.
- [6] D. Dubois and H. Prade. *Possibility Theory*. Plenum Press, New York, NY, USA 1988
- [7] R.A. Fisher. The use of multiple measurements in taxonomic problems. *Annals of Eugenics*, 7(2):179-188. 1936
- [8] Gustafson, E.E. and Kessel, W.C. Fuzzy Clustering with a Fuzzy Covariance Matrix. *IEEE CDC*, San Diego, Californien, pp. 761-766, 1979.
- [9] Höppner, F., Klawonn, F., Kruse, R., and Runkler, T.: *Fuzzy Cluster Analysis*. J. Wiley & Sons, Chichester, England 1999.
- [10] Krishnapuram, R. und Keller, J.: A Possibilistic Approach to Clustering, *IEEE Transactions on Fuzzy Systems*, pp. 98-110, (1) 1993.
- [11] Timm, H.: Fuzzy Cluster Analysis of Classified Data, *IFSA/Nafips 2001*, Vanvouver, to appear.