

# VALIDITY MEASURES FOR HEURISTIC POSSIBILISTIC CLUSTERING

Dmitri A. Viattchenin

United Institute of Informatics Problems of the National Academy of Sciences of Belarus  
Surganov St. 6, 220012 Minsk, Belarus  
e-mail: viattchenin@mail.ru

**Abstract.** A direct algorithm of possibilistic clustering is the effective tool for the data analysis. The approach is based on the concept of allotment among fuzzy clusters. To establish the number of clusters in a data set, some validity measures are presented in this paper. Illustrative examples of application of proposed validity measures to some well-known data sets are given in comparison with a well-known cluster validity index for objective function-based fuzzy clustering algorithms. Preliminary conclusions are formulated.

**Keywords:** clustering, fuzzy cluster, membership degree, allotment, index of fuzziness, density of fuzzy cluster, cluster validity.

## 1. Introduction

The first subsection of this introduction provides a brief review of fuzzy clustering methods. A problem of cluster validity is stated in the second subsection.

### 1.1. Preliminary remarks

Clustering is a process aiming at grouping a set of objects into classes according to the characteristics of data so that objects within a cluster have high mutual similarity while objects in different clusters are dissimilar. Fuzzy sets theory, which was proposed by Zadeh [23], gives an idea of uncertainty of belonging to a cluster, which is described by a membership function. Fuzzy clustering methods have been applied effectively in image processing, data analysis, symbol recognition and modeling. Heuristic methods of fuzzy clustering, hierarchical methods of fuzzy clustering and optimization methods of fuzzy clustering were proposed by different researchers.

The most widespread approach in fuzzy clustering is the optimization approach and the traditional optimization methods of fuzzy clustering are based on the concept of fuzzy  $c$ -partition. Objective function-based fuzzy clustering algorithms can in general be divided into two types: object versus relational. The object data clustering methods can be applied if the objects are represented as points in some multidimensional space. The best-known optimization approach to fuzzy clustering is the method of fuzzy  $c$ -means, developed by Bezdek [2]. The  $FCM$ -algorithm is based on an iterative optimization of the fuzzy objective function, which takes the form:

$$Q_{FCM}(P, \bar{T}) = \sum_{l=1}^c \sum_{i=1}^n u_{li}^\gamma \|x_i - \bar{\tau}^l\|^2, \quad (1)$$

subject to constraints regarding  $u_{li}$

$$\sum_{l=1}^c u_{li} = 1, \quad 0 \leq u_{li} \leq 1, \quad (2)$$

where  $u_{li}$ ,  $l=1, \dots, c$ ,  $i=1, \dots, n$  is the membership degree,  $x_i$ ,  $i \in \{1, \dots, n\}$  is the data point,  $\bar{T} = \{\bar{\tau}^1, \dots, \bar{\tau}^c\}$  is the set fuzzy clusters prototypes and  $\gamma > 1$  is the weighting exponent. Note that the concept of fuzzy  $c$ -partition is defined by the conditions (2). So, the fuzzy  $c$ -partition can be arrayed as a  $(c \times n)$  matrix  $P = [u_{li}]$ .

The  $FCM$ -algorithm is the basis of the family of fuzzy clustering algorithms. These objective function-based fuzzy clustering algorithms were proposed by different authors and they are described by Höppner, Klawonn, Kruse and Runkler [7] in detail.

In the relational approach to fuzzy clustering, the problem of the data classification is solved by expressing a relation which quantifies either similarity, or dissimilarity, between pairs of objects. The most popular examples of fuzzy relational clustering are the Windham's  $AP$ -algorithm [21], the Hathaway, Davenport, and Bezdek's  $RFCM$ -algorithm [6], and the  $ARCA$ -algorithm which was proposed by Corsini, Lazzerini, and Marcelloni in [5]. For example, the  $ARCA$ -algorithm is based on the criterion

$$Q_{ARCA}(P, \bar{T}) = \sum_{l=1}^c \sum_{i=1}^n u_{li}^\gamma \left( \sqrt{\sum_{j=1}^n (d(x_i, x_j) - d(x_j, \bar{\tau}^l))^2} \right)^2, \quad (3)$$

where  $d(x_i, x_j)$  is the dissimilarity relation between the pair of objects  $x_i$  and  $x_j$ , and  $d(x_j, \bar{\tau}^l)$  is the relation between the prototype  $\bar{\tau}^l$ ,  $l \in \{1, \dots, c\}$  and the object  $x_j$ ,  $j \in \{1, \dots, n\}$ .

However, the condition of fuzzy  $c$ -partition is very difficult from essential positions. So, a possibilistic approach to clustering was proposed by Krishnapuram and Keller in [9] and developed by other researchers. This approach can be considered as a way in the optimization approach in fuzzy clustering because all methods of possibilistic clustering are objective function-based methods.

A concept of possibilistic partition is a basis of possibilistic clustering methods and membership values  $\mu_{il}$ ,  $l = 1, \dots, c$ ,  $i = 1, \dots, n$  can be interpreted as the values of typicality degree. For each object  $x_i$ ,  $i = 1, \dots, n$  the grades of membership should satisfy the conditions of a possibilistic partition:

$$\sum_{l=1}^c \mu_{il} > 0, \quad 0 \leq \mu_{il} \leq 1. \quad (4)$$

So, the family of fuzzy sets  $Y(X) = \{A^l \mid l = \overline{1, c}, c \leq n\}$  is the possibilistic partition of the initial set of objects  $X = \{x_1, \dots, x_n\}$  if condition (4) is met. Obviously the conditions of the possibilistic partition (4) are more flexible than the conditions of the fuzzy  $c$ -partition (2).

Heuristic algorithms of fuzzy clustering display high level of essential clarity and low level of a complexity. Some heuristic clustering algorithms are based on a definition of a cluster concept and the aim of these algorithms is cluster detection conform to a given definition. Mandel [10] noted that such algorithms are called algorithms of direct classification or direct clustering algorithms. Direct heuristic algorithms of fuzzy clustering are simple and very effective in many cases. The algorithm of Chiang, Yue, and Yin [4] is a very good illustration for these characterizations.

An outline for a new heuristic method of fuzzy clustering was presented in [13], where concepts of fuzzy  $\alpha$ -cluster and allotment among fuzzy  $\alpha$ -clusters were introduced and a basic version of direct fuzzy clustering algorithm was described. The basic version of direct fuzzy clustering algorithm requires that the number  $c$  of fuzzy  $\alpha$ -clusters be fixed. That is why the basic version of the algorithm, which is described in [13], can be called the  $D-AFC(c)$ -algorithm [17]. Moreover, the allotment of elements of the set of classified objects among fuzzy  $\alpha$ -clusters can be considered as a special case of possibilistic partition (4). These facts were demonstrated in [16] and [18]. So, the  $D-AFC(c)$ -algorithm can be considered as a direct algorithm of possibilistic clustering.

## 1.2. A cluster validity problem

The most important problem of fuzzy clustering is neither the choice of the numerical procedure nor the distance to use but concerns the number  $c$  of fuzzy clusters to look for. Really, lacking in a priori knowledge of the data structure, there is no reason to choose a particular value of  $c$  and one must find a way to measure the acceptance with which cluster structure has been identified by a clustering procedure. This is the so-called cluster validity problem.

The classical approach to cluster validity for fuzzy clustering is based on directly evaluating the fuzzy  $c$ -partition. Measures of cluster validity can be used for the purpose. Many authors have proposed several measures of cluster validity associated with fuzzy  $c$ -partitions. The cluster validity problem can be illustrated by the method of fuzzy  $c$ -means.

Various cluster validity indexes for the  $FCM$ -algorithm were proposed by different researchers [7]. For example, Xie and Beni proposed in [22] well-known validity index, which measures overall average compactness against separation of the fuzzy  $c$ -partition. So, the compactness and separation index is defined in [22] as follows:

$$V_{cs}(P; \bar{T}) = \frac{\sum_{l=1}^c \sum_{i=1}^n u_{il}^{\gamma} \|x_i - \bar{\tau}^l\|^2}{n \times \min_{i \neq l} \|x_i - \bar{\tau}^l\|^2}. \quad (5)$$

The number of clusters that minimizes  $V_{cs}(P; \bar{T})$  is taken as the optimal number  $c$  of fuzzy clusters. The compactness and separation index  $V_{cs}(P; \bar{T})$  is most popular cluster validity criteria for the  $FCM$ -algorithm. Notable that the compactness and separation index  $V_{cs}(P; \bar{T})$  is appropriate for the  $ARCA$ -algorithm, because the  $ARCA$ -algorithm, though being a relational clustering algorithm, generates prototypes. So, the validity measure will be used throughout the paper.

The results of application of the  $D-AFC(c)$ -algorithm to the Anderson's Iris data [1] are considered in [13], [15], [17] and the results show that the  $D-AFC(c)$ -algorithm is a precise and effective numerical procedure for solving classification problems. Moreover, the method of the rapid prototyping fuzzy controllers which is based on deriving fuzzy classification rules from the data on a basis of clustering results obtained from the  $D-AFC(c)$ -algorithm is proposed in [19]. However, validity measures are not proposed for the  $D-AFC(c)$ -algorithm. So, the main goal of this paper is a consideration of the problem of cluster validity for the  $D-AFC(c)$ -algorithm. The contents of this paper is as follows: in the second section, basic concepts of the method are considered and the plan of the  $D-AFC(c)$ -algorithm is proposed, in the third section, methods of evaluation of the fuzzy clusters are considered and cluster

validity indices are proposed, in the fourth section, methods of the data preprocessing are considered and numerical examples of application of the  $D-AFC(c)$  -algorithm with the proposed validity measures to some well-known data sets are given in comparison with the compactness and separation index for objective function-based fuzzy clustering algorithms. In the fifth section some final remarks are stated.

## 2. Outlines of the clustering method

The basic concepts of the heuristic method of possibilistic clustering are considered in the first subsection. A plan of the  $D-AFC(c)$  -algorithm is presented in the second subsection of the section.

### 2.1. Basic concepts

Fuzzy clustering can be considered as a technique of representation of the initial set of objects by fuzzy clusters. The structure of the set of objects can be described by some fuzzy tolerance, that is – a fuzzy binary intransitive relation. So, a fuzzy cluster can be understood as some fuzzy subset originated by fuzzy tolerance relation stipulating that the similarity degree of the fuzzy subset elements is not less than some threshold value. In other words, the value of a membership function of each element of the fuzzy cluster is the degree of similarity of the object to some typical object of fuzzy cluster.

Let us remind the basic concepts of the clustering method based on the concept of allotment among fuzzy clusters, which was proposed in [13]. The concept of fuzzy tolerance is the basis for the concept of fuzzy  $\alpha$ -cluster. That is why definition of fuzzy tolerance must be considered in the first place.

Let  $X = \{x_1, \dots, x_n\}$  be the initial set of elements and  $T: X \times X \rightarrow [0,1]$  some binary fuzzy relation on  $X = \{x_1, \dots, x_n\}$  with  $\mu_T(x_i, x_j) \in [0,1], \forall x_i, x_j \in X$  being its membership function.

**Definition 2.1.** *Fuzzy tolerance is the fuzzy binary intransitive relation which possesses the symmetricity property*

$$\mu_T(x_i, x_j) = \mu_T(x_j, x_i), \forall x_i, x_j \in X, \quad (6)$$

and the reflexivity property

$$\mu_T(x_i, x_i) = 1, \forall x_i \in X. \quad (7)$$

The notions of powerful fuzzy tolerance, feeble fuzzy tolerance and strict feeble fuzzy tolerance were considered in [10], as well. In this context, the classical fuzzy tolerance in the sense of definition 2.1 was called usual fuzzy tolerance. However, the essence of the method here considered does not depend on the kind of fuzzy tolerance. That is why the method herein is described for any fuzzy tolerance  $T$ . A fuzzy tolerance  $T$  on  $X = \{x_1, \dots, x_n\}$  can be represented by a matrix  $T_{n \times n} = [\mu_T(x_i, x_j)], i, j = 1, \dots, n$ .

Let us consider the general definition of fuzzy cluster, the concept of the fuzzy cluster's typical point and the concept of the fuzzy allotment of objects. The number  $c$  of fuzzy clusters can be equal to the number of objects,  $n$ . This is taken into account in further considerations.

Let  $X = \{x_1, \dots, x_n\}$  be the initial set of objects. Let  $T$  be a fuzzy tolerance on  $X$  and  $\alpha$  be  $\alpha$ -level value of  $T, \alpha \in (0,1]$ . Columns or lines of the fuzzy tolerance matrix are fuzzy sets  $\{A^1, \dots, A^n\}$ . Let  $\{A^l, \dots, A^n\}$  be fuzzy sets on  $X$ , which are generated by a fuzzy tolerance  $T$ .

**Definition 2.2.** *The  $\alpha$ -level fuzzy set  $A_{(\alpha)}^l = \{(x_i, \mu_{A^l}(x_i)) \mid \mu_{A^l}(x_i) \geq \alpha, l \in [1, n]\}$  is fuzzy  $\alpha$ -cluster or, simply, fuzzy cluster. So  $A_{(\alpha)}^l \subseteq A^l, \alpha \in (0,1], A^l \in \{A^1, \dots, A^n\}$  and  $\mu_{li}$  is the membership degree of the element  $x_i \in X$  for some fuzzy cluster  $A_{(\alpha)}^l, \alpha \in (0,1], l \in [1, n]$ . Value of  $\alpha$  is the tolerance threshold of fuzzy clusters elements.*

The membership degree of the element  $x_i \in X$  for some fuzzy cluster  $A_{(\alpha)}^l, \alpha \in (0,1], l \in [1, n]$  can be defined as a

$$\mu_{li} = \begin{cases} \mu_{A^l}(x_i), & x_i \in A_{(\alpha)}^l \\ 0, & \text{otherwise} \end{cases}, \quad (8)$$

where an  $\alpha$ -level  $A_{(\alpha)}^l = \{x_i \in X \mid \mu_{A^l}(x_i) \geq \alpha\}, \alpha \in (0,1]$  of a fuzzy set  $A^l$  is the support of the fuzzy cluster  $A_{(\alpha)}^l$ . So, condition  $A_{(\alpha)}^l = \text{Supp}(A_{(\alpha)}^l)$  is met for each fuzzy cluster  $A_{(\alpha)}^l, \alpha \in (0,1], l \in [1, n]$ . Membership degree can be interpreted as a degree of typicality of an element to a fuzzy cluster. The value of a membership function of each element of the fuzzy cluster in the sense of Definition 2.2 is the degree of similarity of the object to some typical object of fuzzy cluster. So, fuzzy clusters in Definition 2.2 are different from fuzzy clusters in the sense (2) from the methodological positions.

In other words, if columns or lines of fuzzy tolerance  $T$  matrix are fuzzy sets  $\{A^1, \dots, A^n\}$  on  $X$ , then fuzzy clusters  $\{A_{(\alpha)}^1, \dots, A_{(\alpha)}^n\}$  are fuzzy subsets of fuzzy sets  $\{A^1, \dots, A^n\}$  for some value  $\alpha, \alpha \in (0,1]$ . The value zero for a fuzzy set membership function is equivalent to non-belonging of an element to a fuzzy set. That is why values of tolerance threshold  $\alpha$  are considered in the interval  $(0, 1]$ .

**Definition 2.3.** *Let  $T$  be a fuzzy tolerance on  $X$ , where  $X$  is the set of elements, and  $\{A_{(\alpha)}^1, \dots, A_{(\alpha)}^n\}$  be the family of fuzzy clusters for some  $\alpha \in (0,1]$ . The point  $\tau_e^l \in A_{(\alpha)}^l$ , for which*

$$\tau_e^l = \arg \max_{x_i} \mu_{li}, \forall x_i \in A_{(\alpha)}^l \quad (9)$$

is called a typical point of the fuzzy cluster  $A_{(\alpha)}^l$ ,  $\alpha \in (0,1]$ ,  $l \in [1, n]$ .

So, the expression (8) defines a possibility distribution function for some  $\tau^l$  over the domain of discourse consisting of all objects  $x_i \in X$ . The distribution will be denoted by  $\pi_l(x_i)$  and the corresponding measure of possibility will be denoted by  $\Pi_l(x_i)$  [18]. Obviously, a typical point of a fuzzy cluster does not depend on the value of tolerance threshold. Moreover, a fuzzy cluster can have several typical points. That is why symbol  $e$  is the index of the typical point.

**Definition 2.4.** *Let*

$R_z^\alpha(X) = \{A_{(\alpha)}^l \mid l = \overline{1, c}, 2 \leq c \leq n, \alpha \in (0,1]\}$  be a family of fuzzy clusters for some value of tolerance threshold  $\alpha$ ,  $\alpha \in (0,1]$ , which are generated by some fuzzy tolerance  $T$  on the initial set of elements  $X = \{x_1, \dots, x_n\}$ . If condition

$$\sum_{l=1}^c \mu_{li} > 0, \quad \forall x_i \in X \tag{10}$$

is met for all fuzzy clusters  $A_{(\alpha)}^l \in R_z^\alpha(X)$ ,  $l = \overline{1, c}$ ,  $c \leq n$ , then the family is the allotment of elements of the set  $X = \{x_1, \dots, x_n\}$  among fuzzy clusters  $\{A_{(\alpha)}^l, l = \overline{1, c}, 2 \leq c \leq n\}$  for some value of the tolerance threshold  $\alpha$ .

It should be noted that several allotments  $R_z^\alpha(X)$  can exist for some tolerance threshold  $\alpha$ . That is why symbol  $z$  is the index of an allotment.

The condition (10) requires that every object  $x_i$ ,  $i = 1, \dots, n$  must be assigned to at least one fuzzy cluster  $A_{(\alpha)}^l$ ,  $l = \overline{1, c}$ ,  $c \leq n$  with the membership degree higher than zero. The condition  $2 \leq c \leq n$  requires that the number of fuzzy clusters in each allotment  $R_z^\alpha(X)$  must be equal or more than two. Otherwise, the unique fuzzy cluster will contain all objects, possibly with different positive membership degrees.

The definition of the allotment among fuzzy clusters (10) is similar to the definition of the possibilistic partition (4). This fact was shown in [16] and [18]. So, the allotment among fuzzy clusters can be considered as the possibilistic partition, and fuzzy clusters in the sense of Definition 2.2 are elements of the possibilistic partition. However, the concept of allotment will be used in further considerations.

The concept of allotment is the central point of the method. But the concept introduced next should be paid attention to, as well.

**Definition 2.5.** *Allotment*

$R_z^\alpha(X) = \{A_{(\alpha)}^l \mid l = \overline{1, n}, \alpha \in (0,1]\}$  of the set of objects among  $n$  fuzzy clusters for some tolerance threshold

$\alpha \in (0,1]$  is the initial allotment of the set  $X = \{x_1, \dots, x_n\}$ .

In other words, if initial data are represented by a matrix of some fuzzy  $T$  then lines or columns of the matrix are fuzzy sets  $A^l \subseteq X$ ,  $l = \overline{1, n}$  and level fuzzy sets  $A_{(\alpha)}^l$ ,  $l = \overline{1, c}$ ,  $\alpha \in (0,1]$  are fuzzy clusters. These fuzzy clusters constitute an initial allotment for some tolerance threshold and they can be considered as clustering components.

Thus, the problem of fuzzy cluster analysis can be defined in general as the problem of discovering the unique allotment  $R^*(X)$ , resulting from the classification process, which corresponds to either most natural allocation of objects among fuzzy clusters or to the researcher's opinion about classification. In the first case, the number of fuzzy clusters  $c$  is not fixed. In the second case, the researcher's opinion determines the kind of the allotment sought and the number of fuzzy clusters  $c$  can be fixed.

If some allotment

$R_z^\alpha(X) = \{A_{(\alpha)}^l \mid l = \overline{1, c}, c \leq n, \alpha \in (0,1]\}$  corresponds to the formulation of a concrete problem, then this allotment is an adequate allotment. In particular, if condition

$$\bigcup_{l=1}^c A_{(\alpha)}^l = X, \tag{11}$$

and condition

$$\text{card}(A_{(\alpha)}^l \cap A_{(\alpha)}^m) = 0, \quad \forall A_{(\alpha)}^l, A_{(\alpha)}^m, l \neq m, \alpha \in (0,1] \tag{12}$$

are met for all fuzzy clusters  $A_{(\alpha)}^l$ ,  $l = \overline{1, c}$  of some allotment  $R_z^\alpha(X) = \{A_{(\alpha)}^l \mid l = \overline{1, c}, c \leq n, \alpha \in (0,1]\}$ , then the allotment is the allotment among fully separate fuzzy clusters.

However, fuzzy clusters in the sense of Definition 2.2 can have an intersection area. This fact was demonstrated in [14]. If the intersection area of any pair of different fuzzy clusters is an empty set, then conditions (11) and (12) are met and fuzzy clusters are called fully separate fuzzy clusters. Otherwise, fuzzy clusters are called particularly separate fuzzy clusters and  $w \in \{0, \dots, n\}$  is the maximum number of elements in the intersection area of different fuzzy clusters. Obviously, for  $w=0$  fuzzy clusters are fully separate fuzzy clusters. So, the conditions (11) and (12) can be generalized for a case of particularly separate fuzzy clusters. Condition

$$\sum_{l=1}^c \text{card}(A_{(\alpha)}^l) \geq \text{card}(X), \quad \forall A_{(\alpha)}^l \in R_z^\alpha(X), \alpha \in (0,1], \text{card}(R_z^\alpha(X)) = c \tag{13}$$

and condition

$$\text{card}(A_{(\alpha)}^l \cap A_{(\alpha)}^m) \leq w, \quad \forall A_{(\alpha)}^l, A_{(\alpha)}^m, l \neq m, \alpha \in (0,1] \tag{14}$$

are generalizations of conditions (11) and (12). Obviously, if  $w=0$  in conditions (13) and (14), then conditions (11) and (12) are met. The adequate allotment  $R_z^\alpha(X)$  for some value of tolerance threshold  $\alpha \in (0,1]$  is a family of fuzzy clusters which are elements of the initial allotment  $R_l^\alpha(X)$  for the value of  $\alpha$  and the family of fuzzy clusters should satisfy the conditions (13) and (14). So, the construction of adequate allotments  $R_z^\alpha(X) = \{A_{(\alpha)}^l | l = \overline{1, c}, c \leq n\}$  for every  $\alpha$  is a trivial problem of combinatorics.

Several adequate allotments can exist. Thus, the problem consists in the selection of the unique adequate allotment  $R^*(X)$  from the set  $B$  of adequate allotments,  $B = \{R_z^\alpha(X)\}$ , which is the class of possible solutions of the concrete classification problem and  $B = \{R_z^\alpha(X)\}$  depends on the parameters of the classification problem. The selection of the unique adequate allotment  $R^*(X)$  from the set  $B = \{R_z^\alpha(X)\}$  of adequate allotments must be made on the basis of evaluation of allotments. The criterion

$$F_1(R_z^\alpha(X), \alpha) = \sum_{l=1}^c \frac{1}{n_l} \sum_{i=1}^{n_l} \mu_{li} - \alpha \cdot c, \quad (15)$$

where  $c$  is the number of fuzzy clusters in the allotment  $R_z^\alpha(X)$  and  $n_l = \text{card}(A_{(\alpha)}^l)$ ,  $A_{(\alpha)}^l \in R_z^\alpha(X)$ , is the number of elements in the support of the fuzzy cluster  $A_{(\alpha)}^l$ , can be used for evaluation of allotments [13]. The criterion

$$F_2(R_z^\alpha(X), \alpha) = \sum_{l=1}^c \sum_{i=1}^{n_l} (\mu_{li} - \alpha), \quad (16)$$

can also be used for evaluation of allotments [16].

Maximum of criterion (15) or criterion (16) corresponds to the best allotment of objects among  $c$  fuzzy clusters. So, the classification problem can be characterized formally as determination of the solution  $R^*(X)$  satisfying

$$R^*(X) = \arg \max_{R_z^\alpha(X) \in B} F(R_z^\alpha(X), \alpha), \quad (17)$$

where  $B = \{R_z^\alpha(X)\}$  is the set of adequate allotments corresponding to the formulation of a concrete classification problem and criteria (15) and (16) are denoted by  $F(R_z^\alpha(X), \alpha)$ .

The criterion (15) can be considered as the average total membership of objects in fuzzy clusters of the allotment  $R_z^\alpha(X)$  minus  $\alpha \cdot c$ . The quantity  $\alpha \cdot c$  regularizes with respect to the number of clusters  $c$  in the allotment  $R_z^\alpha(X)$ . The criterion (16) can be considered as the total membership of objects in fuzzy clusters of the allotment  $R_z^\alpha(X)$  with an appreciation through the value  $\alpha$  of tolerance threshold.

The condition (17) must be met for the some unique allotment  $R_z^\alpha(X) \in B$ . Otherwise, the number  $c$

of fuzzy clusters in the allotment sought  $R^*(X)$  is not appropriate. The important condition was formulated in [14].

## 2.2. A general plan of clustering procedure

The classification problem formulation depends on the parameters of classification and these parameters are determined for a problem of classification in a concrete case. A number  $c$  of fuzzy clusters in the sought allotment  $R^*(X)$  is a unique parameter of the  $D-AFC(c)$ -algorithm. So, the class of possible solutions  $B(c)$  of the classification problem depends on the parameter  $c$  and a unique allotment must be selected from the class  $B(c)$  on a basis of the criteria (15) or (16) calculation for every allotment from the class  $B(c)$ .

There is a seven-step procedure of classification:

1. Calculate  $\alpha$ -level values of the fuzzy tolerance  $T$  and construct the sequence  $0 < \alpha_0 < \alpha_1 < \dots < \alpha_l < \dots < \alpha_z \leq 1$  of  $\alpha$ -levels;
2. Construct the initial allotment  $R_l^\alpha(X) = \{A_{(\alpha)}^l | l = \overline{1, n}\}$ ,  $\alpha = \alpha_l$  for every value  $\alpha_l$  from the sequence  $0 < \alpha_0 < \alpha_1 < \dots < \alpha_l < \dots < \alpha_z \leq 1$ ;
3. Let  $w := 0$ ;
4. Construct allotments  $R_z^\alpha(X) = \{A_{(\alpha)}^l | l = \overline{1, c}, c \leq n\}$ ,  $\alpha = \alpha_l$ , which satisfy conditions (13) and (14) for every value  $\alpha_l$  from the sequence  $0 < \alpha_0 < \alpha_1 < \dots < \alpha_l < \dots < \alpha_z \leq 1$ ;
5. Construct the class of possible solutions of the classification problem  $B(c) = \{R_z^\alpha(X)\}$ ,  $\alpha \in \{\alpha_1, \dots, \alpha_z\}$  for the given number of fuzzy clusters  $c$  and different values of the tolerance threshold  $\alpha$ ,  $\alpha \in \{\alpha_1, \dots, \alpha_z\}$  as follows:
  - if** for some allotment  $R_z^\alpha(X)$ ,  $\alpha \in \{\alpha_1, \dots, \alpha_z\}$  the condition  $\text{card}(R_z^\alpha(X)) = c$  is met
  - then**  $R_z^\alpha(X) \in B(c)$
  - else** let  $w := w + 1$  and go to step 4;
6. Calculate the value of some criterion  $F(R_z^\alpha(X), \alpha)$  for every allotment  $R_z^\alpha(X) \in B(c)$ ;
7. The result  $R^*(X)$  of classification is formed as follows:
  - if** for some unique allotment  $R_z^\alpha(X)$  from the set  $B(c)$  the condition (17) is met
  - then** the allotment is the result of classification
  - else** the number  $c$  of classes is suboptimal.

The allotment  $R_z^\alpha(X) = \{A_{(\alpha)}^l | l = \overline{1, c}, \alpha \in (0,1]\}$  among the given number  $c$  of fuzzy clusters and the corresponding value of tolerance threshold  $\alpha$ ,  $\alpha \in (0,1]$  are the results of classification.

### 3. Cluster validity

Methods for evaluation of the clustering results are given in the first subsection of the section. The second subsection of the section provides three validity measures for the  $D-AFC(c)$ -algorithm.

#### 3.1. Evaluation of the fuzzy clusters

The result of classification must be interpreted from essential positions. Some formal criteria can be useful for the aim. For example, most appropriate distance between fuzzy sets for the data preprocessing can be selected on a basis of the evaluation of the results of classification. A problem of the evaluation of fuzzy clusters was considered in [15].

The qualitative inspection of fuzzy clustering results can be done, e.g., with a linear index of fuzziness or a quadratic index of fuzziness, used for evaluation of fuzziness degree of fuzzy clusters. These two indexes are considered by Kaufmann [8]. So, a modification of the linear index of fuzziness is defined in [15] as

$$I_L(A'_{(\alpha)}) = \frac{2}{n_i} \cdot d_H(A'_{(\alpha)}, \underline{A}'_{(\alpha)}), \quad (18)$$

where  $n_i = \text{card}(A'_{(\alpha)})$ ,  $A'_{(\alpha)} \in R^*(X)$ , is the number of objects in the fuzzy cluster  $A'_{(\alpha)}$  and  $d_H(A'_{(\alpha)}, \underline{A}'_{(\alpha)})$  is the Hamming distance

$$d_H(A'_{(\alpha)}, \underline{A}'_{(\alpha)}) = \sum_{x_i \in A'_i} \left| \mu_{A'_i}(x_i) - \mu_{\underline{A}'_{(\alpha)}}(x_i) \right| \quad (19)$$

between the fuzzy cluster  $A'_{(\alpha)}$  and the crisp set  $\underline{A}'_{(\alpha)}$  nearest to the fuzzy cluster  $A'_{(\alpha)}$ . The membership function of the crisp set  $\underline{A}'_{(\alpha)}$  can be defined as

$$\mu_{\underline{A}'_{(\alpha)}}(x_i) = \begin{cases} 0, & \mu_{A'_i}(x_i) \leq 0.5, \\ 1, & \mu_{A'_i}(x_i) > 0.5, \end{cases} \quad \forall x_i \in A'_i, \quad (20)$$

where  $\alpha \in (0, 1]$ .

The modified quadratic index of fuzziness is defined in [15] as

$$I_Q(A'_{(\alpha)}) = \frac{2}{\sqrt{n_i}} \cdot d_E(A'_{(\alpha)}, \underline{A}'_{(\alpha)}), \quad (21)$$

where  $n_i = \text{card}(A'_{(\alpha)})$ ,  $A'_{(\alpha)} \in R^*(X)$ , and  $d_E(A'_{(\alpha)}, \underline{A}'_{(\alpha)})$  is the Euclidean distance

$$d_E(A'_{(\alpha)}, \underline{A}'_{(\alpha)}) = \sqrt{\sum_{x_i \in A'_i} \left( \mu_{A'_i}(x_i) - \mu_{\underline{A}'_{(\alpha)}}(x_i) \right)^2} \quad (22)$$

between the fuzzy cluster  $A'_{(\alpha)}$  and the crisp set  $\underline{A}'_{(\alpha)}$  which is defined by formula (20).

For each fuzzy cluster  $A'_{(\alpha)}$  in the allotment  $R^*(X)$ , evidently, the following conditions are met:

$$0 \leq I_L(A'_{(\alpha)}) \leq 1, \quad (23)$$

$$0 \leq I_Q(A'_{(\alpha)}) \leq 1. \quad (24)$$

Indexes (18) and (21) show the degree of fuzziness of fuzzy clusters which are elements of the allotment  $R^*(X)$ . Obviously,  $I_L(A'_{(\alpha)}) = I_Q(A'_{(\alpha)}) = 0$  for a crisp set  $A'_{(\alpha)} \in R^*(X)$ . Otherwise, if  $\mu_{A'_i} = 0.5$ ,  $\forall x_i \in A'_i$  then fuzzy cluster  $A'_{(\alpha)} \in R^*(X)$  is a maximally fuzzy set and  $I_L(A'_{(\alpha)}) = I_Q(A'_{(\alpha)}) = 1$ .

The density of fuzzy cluster was defined in [15] as follows:

$$D(A'_{(\alpha)}) = \frac{1}{n_i} \sum_{x_i \in A'_i} \mu_{A'_i}, \quad (25)$$

where  $n_i = \text{card}(A'_{(\alpha)})$ ,  $A'_{(\alpha)} \in R^*(X)$ , and membership degree  $\mu_{A'_i}$  is defined by formula (8). It is obvious that condition

$$0 < D(A'_{(\alpha)}) \leq 1, \quad (26)$$

is met for each fuzzy cluster  $A'_{(\alpha)}$  in  $R^*(X)$ . Moreover,  $D(A'_{(\alpha)}) = 1$  for a crisp set  $A'_{(\alpha)} \in R^*(X)$  for any tolerance threshold  $\alpha$ ,  $\alpha \in (0, 1]$ . The density of fuzzy cluster shows an average membership degree of elements of a fuzzy cluster.

#### 3.2. Validity measures for the D-AFC(c)-algorithm

The most "plausible" number  $c$  of fuzzy clusters in the sought allotment  $R^*(X)$  can be considered as the cluster validity problem for the  $D-AFC(c)$ -algorithm. The number  $c$  of fuzzy clusters and their compactness are contradictory purposes of the classification of  $n$  objects. If compact classes are searched, the most appropriate solution can be obtained with  $n$  classes each consisting of one object. Obviously, such a solution is not useful. So, the number  $c$  of fuzzy clusters must be determined under consideration of the conditions: firstly, the number of fuzzy clusters  $c$  in the sought allotment  $R^*(X)$  must be as possible as less, and, secondly, the membership function of fuzzy clusters of some allotment among  $c$  fuzzy clusters must be sharper than the membership function of fuzzy clusters of allotments for other numbers of fuzzy clusters.

Let  $R_c^*(X)$  be the allotment which corresponds to the result of classification for the given number  $c$  of fuzzy clusters and  $R^c$  be the set of all allotments  $R_c^*(X)$  among  $c$ ,  $c \in \{2, \dots, n\}$  fuzzy clusters. A cluster validity measure can be defined as a mapping  $V: R^c \mapsto \mathfrak{R}$  which can be used to rank the validity of various allotments  $R_c^*(X)$ . Validity measures can be obtained from the indexes which are defined in the previous section.

The fuzziness of the allotment  $R_c^*(X)$  among  $c$  fuzzy clusters can be evaluated as the sum of indexes of fuzziness of fuzzy clusters of the allotment  $R_c^*(X)$ . So, the linear measure of fuzziness of the allotment

must be based on the formula (18), and the measure can be defined as follows:

$$V_{LMF}(R_c^*(X);c) = \sum_{A_{(\alpha)}^l \in R_c^*(X)} \left( I_L(A_{(\alpha)}^l) \right) = \sum_{A_{(\alpha)}^l \in R_c^*(X)} \left( \frac{2}{n_l} \cdot d_H(A_{(\alpha)}^l, \underline{A}_{(\alpha)}^l) \right) \quad (27)$$

The linear measure of fuzziness of the allotment  $V_{LMF}(R_c^*(X);c)$  was proposed in [20].

From other hand, the quadratic measure of fuzziness of the allotment can be defined on the analogy of the linear measure of fuzziness (27):

$$V_{QMF}(R_c^*(X);c) = \sum_{A_{(\alpha)}^l \in R_c^*(X)} \left( I_Q(A_{(\alpha)}^l) \right) = \sum_{A_{(\alpha)}^l \in R_c^*(X)} \left( \frac{2}{\sqrt{n_l}} \cdot d_E(A_{(\alpha)}^l, \underline{A}_{(\alpha)}^l) \right) \quad (28)$$

where  $I_Q(A_{(\alpha)}^l)$  is the modified quadratic index of fuzziness (21).

The justification of both measures of fuzziness is intuitive. It is obvious that the fuzziness of each fuzzy cluster  $A_{(\alpha)}^l \in R_c^*(X)$  depends on the size of the fuzzy cluster. The number of objects  $n_l$  in each fuzzy cluster  $A_{(\alpha)}^l \in R_c^*(X)$  is decreasing with increasing of the number  $c$  of fuzzy clusters in the allotment  $R_c^*(X)$ . That is why the fuzziness of each fuzzy cluster  $A_{(\alpha)}^l \in R_c^*(X)$  is decreasing with increasing of the number  $c$  of fuzzy clusters. In other words, for  $c \rightarrow n$  we have  $n_l \rightarrow 1$  and  $I_L(A_{(\alpha)}^l) \rightarrow 0$ ,  $I_Q(A_{(\alpha)}^l) \rightarrow 0$  for all  $A_{(\alpha)}^l$ ,  $l \in \{1, \dots, c\}$ . So, for  $c \rightarrow n$  we have  $V_{LMF}(R_c^*(X);c) \rightarrow 0$  and  $V_{QMF}(R_c^*(X);c) \rightarrow 0$ . So, the maximal value of a measure of fuzziness of the allotment  $R_c^*(X)$  corresponds to the minimal number  $c$  of compact fuzzy clusters  $A_{(\alpha)}^l \in R_c^*(X)$ ,  $l=1, \dots, c$  in the sought allotment. Using  $V_{LMF}(R_c^*(X);c)$  or  $V_{QMF}(R_c^*(X);c)$ , the optimal number of fuzzy clusters can be obtained by maximizing the index value.

The density of fuzzy cluster (25) can be considered as the basis for a validity measure, too. The validity measure must take into account the compactness of fuzzy clusters which is characterized by their density. The density of each fuzzy cluster  $A_{(\alpha)}^l \in R_c^*(X)$  is increasing with increasing of the number  $c$  of fuzzy clusters. So, for  $c \rightarrow n$  we have  $D(A_{(\alpha)}^l) \rightarrow 1$  for all  $A_{(\alpha)}^l$ ,  $l \in \{1, \dots, c\}$  and for  $c \rightarrow n$  we have  $\alpha \rightarrow 1$ . Thus, the value of the tolerance threshold  $\alpha$  must be taken into account. So, the validity measure can be defined as the ratio of the sum of densities of fuzzy clusters of some allotment to the number of fuzzy clusters minus the value of the tolerance threshold  $\alpha$ . However, a case of particularly separate fuzzy clusters must be

taken into account. That is why the sum of membership degrees of elements in intersection areas of fuzzy clusters must be calculated and the ratio of the number  $c$  of fuzzy clusters in the allotment  $R_c^*(X)$  to the number of elements of the data set must be taken into account, too. So, the measure of separation and compactness of the allotment can be defined in the following way:

$$V_{MSC}(R_c^*(X);c) = \frac{\sum_{A_{(\alpha)}^l \in R_c^*(X)} D(A_{(\alpha)}^l)}{c} + \frac{c}{n} \sum_{x_j \in \Theta} \mu_{x_j} - \alpha \quad (29)$$

where  $\Theta$  is a set of elements  $x_j$ ,  $j \in \{1, \dots, n\}$  in all intersection areas of different fuzzy clusters.

Thus, the validity measure (29) has three components. The first component is the ratio of the sum of densities of fuzzy clusters of an allotment  $R_c^*(X)$  to the number of fuzzy clusters  $c$ , the second is a penalty term which regularizes with respect to membership values of elements in intersection areas of fuzzy clusters, and the third component is included to make the value of the tolerance threshold. Note that the value of the second component of the measure of separation and compactness (29) be equal to zero in the case of fully separate fuzzy clusters in the allotment  $R_c^*(X)$ . From other hand, several allotments  $R_c^*(X) \in B(c)$  among  $c$  fuzzy clusters cannot exist as solutions of the classification problem because the condition (17) must be met for the some unique allotment  $R_c^*(X) \in B(c)$ . The measure of separation and compactness of the allotment  $V_{MSC}(R_c^*(X);c)$  increases when  $c$  is closer to  $n$ . Thus the optimum value of  $c$  is obtained by minimizing  $V_{MSC}(R_c^*(X);c)$  over  $c = 2, \dots, c_{\max}$  where  $2 < c_{\max} < n$ .

## 4. Experimental results

Results of some well-known data sets processed by the  $D-AFC(c)$ -algorithm using the proposed validity measures in comparison with two objective-function fuzzy clustering algorithms and corresponding validity measures are presented in the section. The first subsection includes the Tamura's relational data description and results of their processing by the  $D-AFC(c)$ -algorithm using the proposed validity measures in comparison with the  $ARCA$ -algorithm. In the second subsection the Anderson's Iris data is used for testing of proposed validity measures in comparison with the  $FCM$ -algorithm.

### 4.1. The Tamura's portrait data

Let us consider an application of proposed validity measures to the classification problem for the following illustrative example. The problem of

classification of family portraits coming from three families was considered by Tamura, Higuchi and Tanaka in [12]. The number of portraits was equal to 16 and the real portrait assignment among three classes is presented in Figure 1.

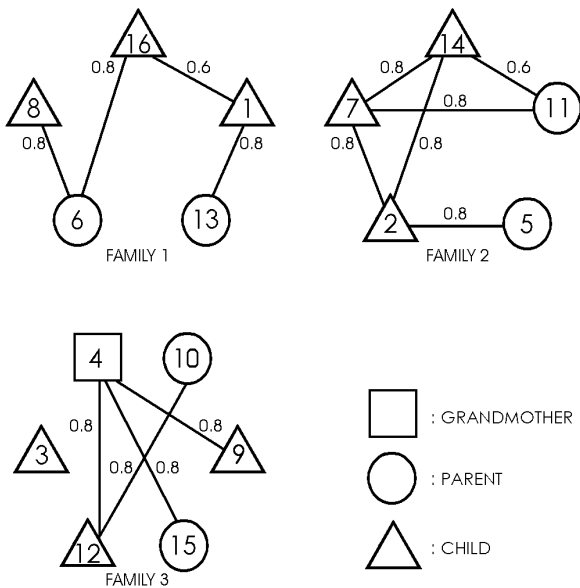


Figure 1. Real portraits classification

Table 1. The matrix of subjective similarities

$i$	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16
1	1.0															
2	0.0	1.0														
3	0.0	0.0	1.0													
4	0.0	0.0	0.4	1.0												
5	0.0	0.8	0.0	0.0	1.0											
6	0.5	0.0	0.2	0.2	0.0	1.0										
7	0.0	0.8	0.0	0.0	0.4	0.0	1.0									
8	0.4	0.2	0.2	0.5	0.0	0.8	0.0	1.0								
9	0.0	0.4	0.0	0.8	0.4	0.2	0.4	0.0	1.0							
10	0.0	0.0	0.2	0.2	0.0	0.0	0.2	0.0	0.2	1.0						
11	0.0	0.5	0.2	0.2	0.0	0.0	0.8	0.0	0.4	0.2	1.0					
12	0.0	0.0	0.2	0.8	0.0	0.0	0.0	0.0	0.4	0.8	0.0	1.0				
13	0.8	0.0	0.2	0.4	0.0	0.4	0.0	0.4	0.0	0.0	0.0	0.0	1.0			
14	0.0	0.8	0.0	0.2	0.4	0.0	0.8	0.0	0.2	0.2	0.6	0.0	0.0	1.0		
15	0.0	0.0	0.4	0.8	0.0	0.2	0.0	0.0	0.2	0.0	0.0	0.2	0.2	0.0	1.0	
16	0.6	0.0	0.0	0.2	0.2	0.8	0.0	0.4	0.0	0.0	0.0	0.0	0.4	0.2	0.0	1.0

The  $D-AFC(c)$ -algorithm was applied to the matrix of fuzzy tolerance for  $c=2, \dots, 5$  using the proposed validity measures. The performance of the proposed validity measures is shown in Figures 2 – 4.

The actual number of fuzzy clusters is equal 3 and this number corresponds to the maximum of the linear measure of fuzziness of the allotment  $V_{LMF}(R_c^*(X);c)$  and the maximum of the quadratic measure of fuzziness of the allotment  $V_{QMF}(R_c^*(X);c)$ . From other hand,

The matrix of subjective similarities contains the results of a study aimed at discovering the similarity degree of 16 people belonging to three families. The study was conducted by showing the photo of the 16 members of three families to experts who did not know them. The subjective similarities assigned to the individual pairs of portraits collected in the tabular format are presented in Table 1, where  $x_i, i=1, \dots, 16$  identify the 16 people to be grouped into families. In fact, the matrix of subjective similarities is the matrix of a fuzzy tolerance and the  $D-AFC(c)$ -algorithm can be applied to the matrix directly. Obviously, the matrix which is presented in Table 1 has no metric characteristic.

The data were originally analyzed in order to identify families with the technique of first transforming the matrix of a fuzzy tolerance into a matrix of a fuzzy similarity relation and then taking an appropriate  $\alpha$ -cut of the fuzzy similarity relation [12]. The best partition proved to be obtained with  $\alpha$ -cut equal to 0.6. The partition identified the following three families  $A^1 = \{x_1, x_6, x_8, x_{13}, x_{16}\}$ ,  $A^2 = \{x_2, x_5, x_7, x_{11}, x_{14}\}$  and  $A^3 = \{x_4, x_9, x_{10}, x_{12}, x_{15}\}$ . However, person  $x_3$  is not a member of any of the three families.

the minimal value of the measure of separation and compactness of the allotment is equal to 0.4385 and this value corresponds to the four fully separate fuzzy clusters. The value of the total number of elements in intersection areas is equal to 3 for the allotment among two particularly separate fuzzy clusters, for  $c = 3$  we have the total number of elements in intersection areas 2 and the value of the total number of elements in intersection areas is equal to 1 for  $c = 5$ .

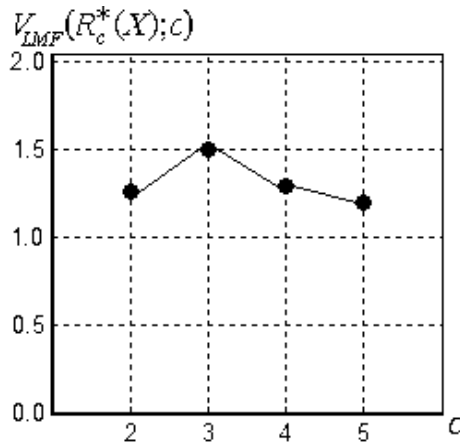


Figure 2. Plot of the linear measure of fuzziness of the allotment as a function of the number of clusters

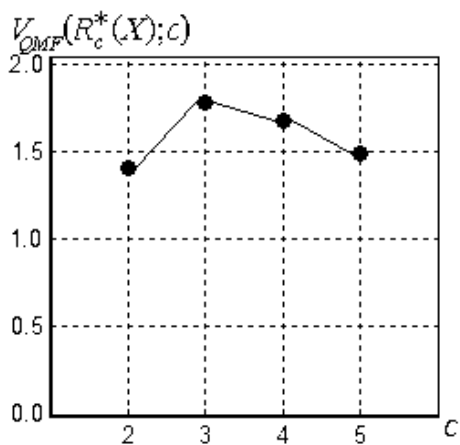


Figure 3. Plot of the quadratic measure of fuzziness of the allotment as a function of the number of clusters

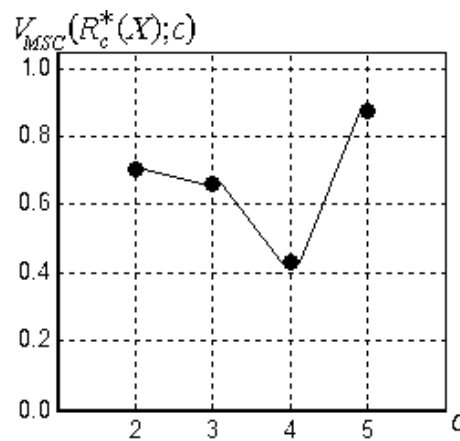


Figure 4. Plot of the measure of separation and compactness of the allotment as a function of the number of clusters

The application of the  $D-AFC(c)$ -algorithm to the classification problem was made in comparison with the  $ARCA$ -algorithm of fuzzy clustering [5] using the compactness and separation index (5) for  $c=2, \dots, 5$ . In order to compare proposed validity measures with the relational  $ARCA$ -algorithm of fuzzy clustering, we

transformed the initial matrix into a dissimilarity matrix by complementing the relationship degrees. The performance of the compactness and separation index is shown in Figure 5.

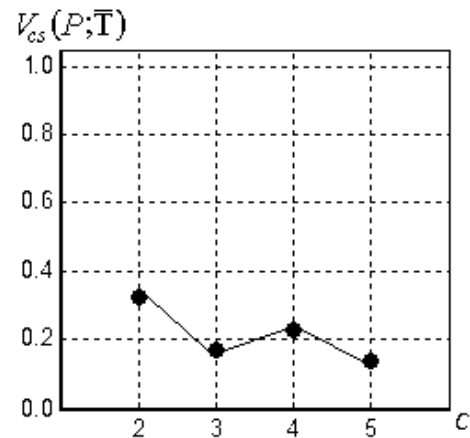


Figure 5. Plot of the compactness and separation index as a function of the number of clusters

We observed that the minimal value of the compactness and separation index corresponds to the five fuzzy clusters. The optimal number of fuzzy clusters is equal to 3 and this number corresponds to the first minimum of the compactness and separation index  $V_{cs}(P; \bar{T})$ .

#### 4.2. The Anderson's Iris data

The Anderson's [1] Iris data is the most known database to be found in the pattern recognition literature. The data set represents different categories of Iris plants having four attribute values. The four attribute values represent the sepal length, sepal width, petal length and petal width measured for 150 irises. It has three classes Setosa, Versicolor and Virginica, with 50 samples per class. The problem is to classify the plants into three subspecies on the basis of this information. It is known that two classes Versicolor and Virginica have some amount of overlap while the class Setosa is linearly separable from the other two.

Let us consider the effectiveness of the proposed clustering validity measures by testing the Anderson's Iris data set. However, a method for the data preprocessing for using the  $D-AFC(c)$ -algorithm must be considered [17].

The Anderson's Iris data can be presented as a matrix of attributes  $\hat{X}_{150 \times 4} = [\hat{x}_i^t]$ ,  $i = 1, \dots, 150$ ,  $t = 1, \dots, 4$ , where the value  $\hat{x}_i^t$  is the value of the  $t$ -th attribute for  $i$ -th object. The data can be normalized as follows:

$$x_i^t = \frac{\hat{x}_i^t}{\max_i \hat{x}_i^t}, \quad (30)$$

for all attributes  $x^t$ ,  $t = 1, \dots, m$ . So, each object can be considered as a fuzzy set  $x_i$ ,  $i = 1, \dots, n$ , and

$x_i^t = \mu_{x_i}(x^t) \in [0,1]$ ,  $i=1,\dots,n$ ,  $t=1,\dots,m$  are their membership functions. After application of the squared normalized Euclidean distance [8]

$$\varepsilon(x_i, x_j) = \frac{1}{m} \sum_{t=1}^m (\mu_{x_i}(x^t) - \mu_{x_j}(x^t))^2, \quad i, j = \overline{1, n}, \quad (31)$$

to the matrix of normalized data  $X_{n \times m} = [\mu_{x_i}(x^t)]$ ,  $i=1,\dots,n$ ,  $t=1,\dots,m$  a matrix of a fuzzy intolerance  $I = [\mu_I(x_i, x_j)]$ ,  $i, j = 1,\dots,n$  is obtained. The matrix of fuzzy tolerance  $T = [\mu_T(x_i, x_j)]$ ,  $i, j = 1,\dots,n$  is obtained after application of complement operation

$$\mu_T(x_i, x_j) = 1 - \mu_I(x_i, x_j), \quad \forall i, j = 1,\dots,n, \quad (32)$$

to the matrix of fuzzy intolerance  $I = [\mu_I(x_i, x_j)]$ ,  $i, j = 1,\dots,n$  obtained from previous operations.

We applied the  $D-AFC(c)$ -algorithm to the obtained matrix of fuzzy tolerance for  $c=2,\dots,5$ . So, we calculated the values of the proposed validity measures for different  $c$  values and we plotted these validity measures in Figures 6 – 8.

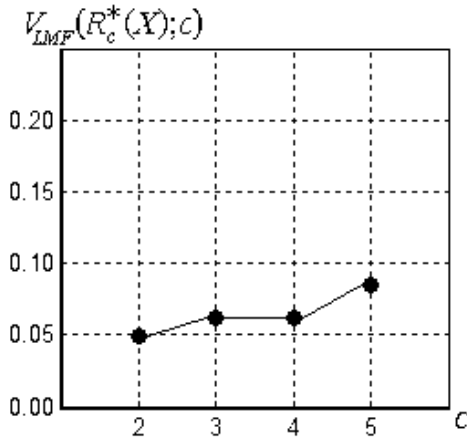


Figure 6. Plot of the linear measure of fuzziness of the allotment as a function of the number of clusters

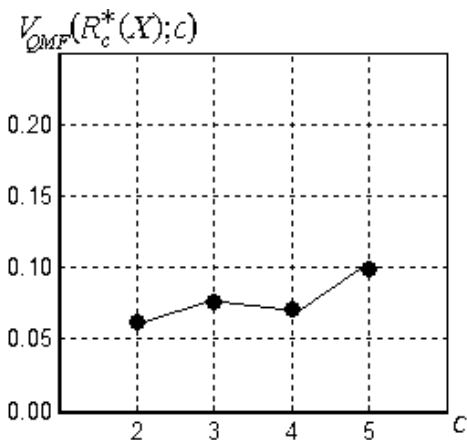


Figure 7. Plot of the quadratic measure of fuzziness of the allotment as a function of the number of clusters

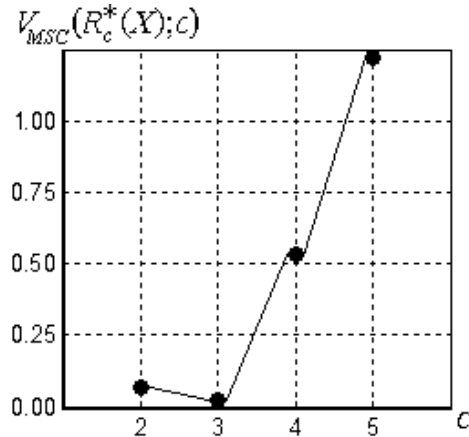


Figure 8. Plot of the measure of separation and compactness of the allotment as a function of the number of clusters

By executing the  $D-AFC(c)$ -algorithm for  $c=2,\dots,5$ , we obtain that the optimal cluster number  $c$  is chosen at  $c=5$  for the linear measure of fuzziness of the allotment and the quadratic measure of fuzziness of the allotment. However, the number of fuzzy clusters  $c=3$  corresponds to the first maximum for both validity measures. From other hand, the measure of separation and compactness of the allotment finds the optimal cluster number  $c$  at  $c=3$ . Allotments among fully separated fuzzy clusters were obtained for  $c=2$  and  $c=3$ . The value of the total number of elements in intersection areas is equal to 10 for the allotment among four particularly separate fuzzy clusters and the value of the total number of elements in intersection areas is equal to 18 for  $c=5$ . So, the results of proposed validity measures seem to be appropriate.

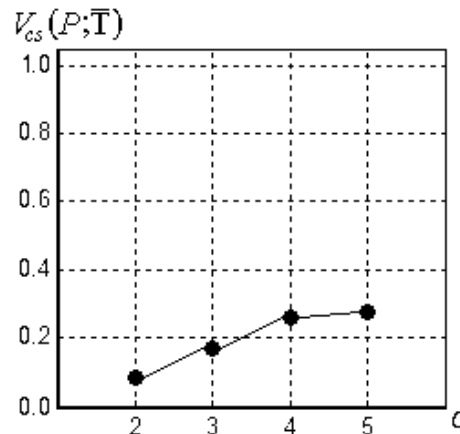


Figure 9. Plot of the compactness and separation index as a function of the number of clusters

For the comparison purpose, the Iris data set was also tested by the  $FCM$ -algorithm using the compactness and separation index (5). The  $FCM$ -algorithm was applied to the data set with the weighting exponent  $\gamma=2.0$  and the value of a small threshold  $\varepsilon=0.001$  for  $c=2,\dots,5$ . The performance of

the compactness and separation index for the Iris data set is shown in Figure 9.

The optimal cluster number  $c$  is chosen at  $c=2$  for the compactness and separation index. Note that most validity measures reported in the literature provide two clusters for these data [3].

## 5. Concluding remarks

Results of experiments are summarized and discussed in the first subsection of the section. The second subsection deals with the perspectives on future investigations.

### 5.1. Discussion

In conclusion, it should be said that fuzzy cluster and allotment concepts have an epistemological motivation. That is why the results of application of the heuristic possibilistic clustering method based on the allotment concept can be very well interpreted. The  $D-AFC(c)$ -algorithm can be applied directly to the data given as the matrix of tolerance coefficients. This means that it can be used with the object by attributes data, by choosing a suitable metric to measure similarity or it can be used in situations where object by object proximity data are available. Moreover, the  $D-AFC(c)$ -algorithm depends on the set of adequate allotments only. That is why the clustering results are stable.

Cluster validity measures are introduced in the paper and numerical experiments confirmed their utility. Some well-known data sets are used for illustrating the properties of the  $D-AFC(c)$ -algorithm. The proposed validity measures provide useful information about structure of the data. Thus, the results of application of the proposed validity measures to the data sets show that these validity measures are effective tools for solving the classification problem. However, the behavior of proposed cluster validity measures has not been justified from mathematical positions. Moreover, it is impossible to judge which one is the best one and the most appropriate with respect to the number of fuzzy clusters. Each of the mentioned measures works well with a certain class of data in common. So, we can conclude that the use of some one validity measure may produce serious hesitation. It will be a reasonable way to make use of various validity measures and compare the obtained clustering results.

### 5.2. Perspectives

Pedrycz [11] noted that the behavior of validity measures has not been theoretically justified, but simulation experiments confirmed their utility. So, for the most appropriate number of fuzzy clusters in the allotment an extremal value of an index or a significant jump of its values can be observed. Moreover, the

comparison of the results, obtained from various fuzzy clustering methods with their cluster validity indexes, can be considered as a useful approach to the data analysis.

Some other validity measures can be proposed for the  $D-AFC(c)$ -algorithm. In the first place, the value of the total number of elements in intersection areas can be introduced in the linear measure of fuzziness of the allotment and the quadratic measure of fuzziness of the allotment. So, modifications of the corresponding validity measures can be proposed. In the second place, the results of the application of the  $D-AFC(c)$ -algorithm to the data depend on the selected distance [17].

These perspectives for investigations are of great interest both from the theoretical point of view and from the practical one as well.

### Acknowledgements

The investigations were made in the Ostfalia University of Applied Sciences under the German Academic Exchange Service financial support. I am grateful to Prof. Frank Klawonn for his interest in my investigations and useful remarks during the paper preparation. I also thank Mr. Aliaksandr Damaratski for elaborating experimental software. I would like to thank the anonymous referees for their valuable comments.

### References

- [1] E. Anderson. The irises of the Gaspé Peninsula. *Bulletin of the American Iris Society*, 1935, Vol.59, 2-5.
- [2] J.C. Bezdek. Pattern Recognition with Fuzzy Objective Function Algorithms. *Plenum Press, New York*, 1981.
- [3] M. Bouguessa, S. Wang, H. Sun. An objective approach to cluster validation. *Pattern Recognition Letters*, 2006, Vol.27, No.13, 1419-1430.
- [4] J.-H. Chiang, S. Yue, Z.-X. Yin. A new fuzzy cover approach to clustering. *IEEE Transactions on Fuzzy Systems*, 2004, Vol.12, No.2, 199-208.
- [5] P. Corsini, B. Lazzerini, F. Marcelloni. A new fuzzy relational clustering algorithm based on the fuzzy C-means algorithm. *Soft Computing*, 2005, Vol.9, No.6, 439-447.
- [6] R.J. Hathaway, J.W. Davenport, J.C. Bezdek. Relational duals of the C-means clustering algorithms. *Pattern Recognition*, 1989, Vol.22, No.2, 205-212.
- [7] F. Höppner, F. Klawonn, R. Kruse, T. Runkler. Fuzzy Cluster Analysis: Methods for Classification, Data Analysis and Image Recognition. *Wiley Inter-science, Chichester*, 1999.
- [8] A. Kaufmann. Introduction to the Theory of Fuzzy Subsets. *Academic Press, New York*, 1975.
- [9] R. Krishnapuram, J.M. Keller. A possibilistic approach to clustering. *IEEE Transactions on Fuzzy Systems*, 1993, Vol.1, No.2, 98-110.
- [10] I.D. Mandel. Clustering Analysis. *Finansy i Statistika, Moscow*, 1988 (in Russian).

- [11] **W. Pedrycz.** Fuzzy sets in pattern recognition: methodology and methods. *Pattern Recognition*, 1990, Vol.23, No.1/2, 121-146.
- [12] **S. Tamura, S. Higuchi, K. Tanaka.** Pattern classification based on fuzzy relations. *IEEE Transactions on Systems, Man, and Cybernetics*, 1971, Vol.1, No.1, 61-66.
- [13] **D.A. Viattchenin.** A new heuristic algorithm of fuzzy clustering. *Control and Cybernetics*, 2004, Vol.33, No.2, 323-340.
- [14] **D.A. Viattchenin.** On the number of fuzzy clusters in the allotment. *Proceedings of the 7<sup>th</sup> International Conference on Computer Data Analysis and Modeling (CDAM'2004)*, Minsk, Belarus, 2004, Vol.1, 198-201.
- [15] **D.A. Viattchenin.** On the inspection of classification results in the fuzzy clustering method based on the allotment concept. *Proceedings of the 4<sup>th</sup> International Conference on Neural Networks and Artificial Intelligence (ICNNAI'2006)*, Brest, Belarus, 2006, 210-216.
- [16] **D.A. Viattchenin.** A direct algorithm of possibilistic clustering with partial supervision. *Journal of Automation, Mobile Robotics and Intelligent Systems*, 2007, Vol.1, No.3, 29-38.
- [17] **D.A. Viattchenin.** A methodology of fuzzy clustering with partial supervision. *Systems Science*, 2007, Vol.33, No.4, 61-71.
- [18] **D.A. Viattchenin.** On possibilistic interpretation of membership values in fuzzy clustering method based on the allotment concept. *Proceedings of the Institute of Modern Knowledge*, 2008, No.3, 85-90 (in Russian).
- [19] **D.A. Viattchenin.** Outlines for a new approach to generating fuzzy classification rules through clustering techniques. *Proceedings of the 10<sup>th</sup> International Conference on Pattern Recognition and Information Processing (PRIP'2009)*, Minsk, Belarus, 2009, 82-87.
- [20] **D.A. Viattchenin, A. Damaratski, D. Novikau.** Relational clustering of heterogeneous fuzzy data. *D.A. Viattchenin (ed.), Developments in Fuzzy Clustering*, VEVER Publishing House, Minsk, 2009, 76-91.
- [21] **M.P. Windham.** Numerical classification of proximity data with assignment measures. *Journal of Classification*, 1985, Vol.2, No.1, 157-172.
- [22] **X.L. Xie, G.A. Beni.** A validity measure for fuzzy clustering. *IEEE Transactions on Pattern Analysis and Machines Intelligence*, 1991, Vol.13, No.8, 841-847.
- [23] **L.A. Zadeh.** Fuzzy sets. *Information and Control*, 1965, Vol.8, No.3, 338-353.

Received February 2010.

DOI: 10.5755/j01.itc.39.4.12390