

A Cluster Validity Index for Fuzzy Clustering

Soon H. Kwon, Haiyoung Lee and Ick Choy*

School of Electrical and Electronic Eng., Yeungnam University

**Korea Institute of Science and Technology*

ABSTRACT

In this paper, a new cluster validation index which is heuristic but able to eliminate the monotonically decreasing tendency occurring in which the number of cluster c gets very large and close to the number of data points n is proposed. We review the FCM algorithm and some conventional cluster validity criteria, discuss on the limiting behavior of the proposed validity index, and provide some numerical examples showing the effectiveness of the proposed cluster validity index.

1. Introduction

Since Zadeh's formulation of fuzzy set theory, many fuzzy set-based approaches to fields such as control, pattern recognition, decision making, and clustering have been developed and applied to systems with uncertainty. The basic idea of these approaches is to represent the uncertainty of the given systems by means of fuzzy rules and their membership functions defined over appropriate discourses. One of the most prominent applications of it may be a fuzzy logic-based modeling by means of fuzzy clustering [10].

Cluster analysis is to place elements into groups or clusters suggested by a given data set $X = \{x_1, \dots, x_n\} \subset R^p$ which are n points in the p -dimensional space for summarizing data or finding "natural" or "real" substructures in the data set. The Fuzzy C-Means (FCM) algorithm [1] and its derivatives based on the possibilistic approach [3,4] for the cluster analysis have been the dominant approaches in both theory and practical applications of fuzzy techniques to unsupervised classification for the last two decades.

As pointed out in the literature of Milligan [2], a cluster analysis will not only refer to clustering methods such as the FCM and the possibilistic approach but also to the overall sequence of steps such as clustering elements, clustering variables, variable standardization, measure of association, number of clusters, interpretation, testing, and replication. In recent years, many literatures have paid a great deal of attention to cluster validity issues and many functionals have been proposed for validation of partitions of data produced by the FCM algorithm [5-9]. According to the

Pal and Bezdek's analysis [9], the Fukuyama-Sugeno index [6] is sensitive to both high and low values of the weighting exponent m and may be unreliable because of this. The Xie-Beni index provided the best response over a wide range of choices for the number of clusters, (2-10), and for the weighting exponent m from 1.01-7. On the basis of their analysis, they suggested that the best choice for the weighting exponent m may be probably in the interval [1.5, 2.5], whose mean and midpoint, $m = 2$, have often been the preferred choice for many users of the FCM.

However, the Xie-Beni index v_{XB} has a flaw that is monotonically decreasing when the number of cluster c gets very large and close to the number of data points n . Xie and Beni suggested that an ad hoc punishing function should be imposed to eliminate the monotonically decreasing tendency, but not discussed how to choose the function. It is highly recommended to impose a punishing function, which is a function of the number of cluster, to eliminate the monotonically decreasing tendency of the cluster validation indexes as like as the statistical model selection criteria do.

In this paper, we propose a new cluster validity index, which is an extension of the Xie-Beni index v_{XB} , with a term of a punishing function to overcome the problem of the monotonically decreasing tendency of the conventional cluster validity indexes. The proposed index has a property to be able to eliminate the decreasing tendency occurring in which the number of cluster c gets very large and close to the number of data points n .

This paper is organized as follows. In section 2, we review the FCM algorithm and some cluster validity

criteria, and present a new cluster validity index. Section 3 describes some numerical examples showing the effectiveness of the presented cluster validity measure.

2. Fuzzy c-means algorithm and cluster validity

The FCM algorithm is a constrained optimization problem which minimizes the following objective function with respect to membership functions u_{ij} and cluster centroid v_i ,

$$J_m(U, V; X) = \sum_{j=1}^n \sum_{i=1}^c u_{ij}^m \|x_j - v_i\|^2 \quad (1)$$

where $U = [u_{ij}]$ is a $c \times n$ matrix, c is the number of clusters, n is the number of data points, satisfying the conditions in (2),

$$M_{fcn} = \left\{ U \in R^{cn} \mid u_{ij} \in [0, 1] \forall i, j; 0 < \sum_{j=1}^n u_{ij} < n \forall i, \text{ and } \sum u_{ij} = 1 \forall j \right\} \quad (2)$$

$V = (v_1, \dots, v_c)$ is a vector of cluster centers, $v_i \in R^p$ for $c \geq i \geq 1$ and $\|\bullet\|$ denotes any inner product norm. Optimal partitions U^* of X are taken from pairs (U^*, V^*) that are local minimizers of J_m obtained by iteration through the following necessary conditions.

Fuzzy c-means theorem [1]: If

$$D_{ij} = \|x_j - v_i\|_A > 0 \quad \forall i, j,$$

the weighting exponent $m > 1$, and a data set X contains $c < n$ distinct points, then $(U, V) \in M_{fcn} \times R^{cp}$ may minimize J_m only if

$$u_{ij} = \left[\sum_{k=1}^c \left(\frac{D_{ijA}}{D_{jkA}} \right)^{\frac{2}{m-1}} \right]^{-1}, \quad 1 \leq i \leq c; 1 \leq j \leq n$$

$$v_i = \frac{\sum_{j=1}^n u_{ij}^m x_j}{\sum_{j=1}^n u_{ij}^m}, \quad 1 \leq i \leq c. \quad (3)$$

If for some i and j , $D_{ijA} = 0$, a singularity occurs, then assign 0's to each u_{ij} for which $D_{ijA} > 0$, and distribute membership functions arbitrary across the x_k 's for which $D_{ijA} = 0$, subject to the constraints in (2). Some limiting properties of (3) have been studied by Pal and Bezdek

[9] and are not discussed here.

2.1 Conventional cluster validity indexes

Among a class of cluster validity functionals such as the Dunn's normalized partition entropy [5] :

$$v_D(U) = \frac{n}{n-c} v_{PE} = -\frac{1}{n-c} \sum_{j=1}^n \sum_{i=1}^c u_{ij} \log_a(u_{ij}), \quad (4)$$

the Bezdek's partition coefficient [1] :

$$v_{PC}(U) = \frac{1}{n} \sum_{j=1}^n \sum_{i=1}^c u_{ij}^2 \quad (5)$$

the Bezdek's partition entropy [1] :

$$v_{PE}(U) = -\frac{1}{n} \sum_{j=1}^n \sum_{i=1}^c u_{ij} \log_a(u_{ij}), \quad (6)$$

where logarithmic base $a \in (1, \infty)$ and $u_{ij} \log(u_{ij}) \cong 0$ whenever $u_{ij} = 0$,

the Fukuyama-Sugeno index [6] :

$$v_{FS}(U, V; X) = \sum_{j=1}^n \sum_{i=1}^c u_{ij}^m (\|x_j - v_i\|^2 - \|v_i - \bar{v}\|^2), \quad (7)$$

the Xie-Beni index [8]:

$$v_{XB}(U, V; X) = \frac{\sum_{j=1}^n \sum_{i=1}^c u_{ij}^m \|x_j - v_i\|^2}{n \left[\min_{i \neq k} (\|v_i - v_k\|^2) \right]}, \quad (8)$$

and the extended FCM Xie-Beni index [8] :

$$v_{XB}(U, V; X) = \frac{\sum_{j=1}^n \sum_{i=1}^c u_{ij}^m \|x_j - v_i\|^2}{n \left[\min_{i \neq k} (\|v_i - v_k\|^2) \right]}, \quad (9)$$

we consider only the Xie-Beni index v_{XB} given by (8) because it provides the best response over a wide range of choices for the number of clusters and for the weighting exponent m as discussed in the introduction.

Xie and Beni stated that v_{XB} decreases monotonically when the number of clusters c is close to n . To avoid the indetermination due to the monotonicity, they recommended plotting v_{XB} as a function of c , finding the starting point of the monotonic epoch as the maximum cluster number to be considered, and then selecting a value c minimizing v_{XB} . Because it requires a cumbersome procedure to find an optimum value of c , it may not be a sufficiently good cluster validity index even though it provides good responses through the cumbersome procedures.

2.2 Extension of the Xie-Beni index

In clustering, it attempts to maximize intra-class similarity and inter-class differences. In this sense, a new cluster validity index v_k is defined as

$$v_k(U, V; X) = \frac{\sum_{j=1}^n \sum_{i=1}^c u_{ij}^2 \|x_j - v_i\|^2 + \frac{1}{c} \sum_{i=1}^c \|v_i - \bar{v}\|^2}{\min_{i \neq k} (\|v_i - v_k\|^2)}, \quad (10)$$

$$\text{where } \bar{v} = \frac{1}{n} \sum_{j=1}^n x_j.$$

The first term of the numerator in (10) measures the intra-class similarity, that is, how compact each and every class is. The more similar (compact) the classes, the smaller it is. It is independent of the number of data points. The second term of the numerator in (10) is an ad hoc punishing function imposed to eliminate the decreasing tendency occurring when the number of cluster c gets very large and close to the number of data points n . The denominator in (10) which is the minimum distance between cluster centroids measures the inter-class difference. A larger value of it indicates that every cluster is well-separated. Our goal is to find the fuzzy c -partition with the smallest value of v_k .

In order to investigate the limiting behavior of the proposed index, we take a limit of the validity functional as c approaches n .

Xie-Beni index, $c \rightarrow n$: Since

$$\lim_{c \rightarrow n} \|x_j - v_i\|^2 = 0, \quad (11)$$

we have

$$\lim_{c \rightarrow n} v_{XB}(U, V; X) = \lim_{c \rightarrow n} \frac{\sum_{j=1}^n \sum_{i=1}^c u_{ij}^2 \|x_j - v_i\|^2}{n \left[\min_{i \neq k} (\|v_i - v_k\|^2) \right]} = 0 \quad (12)$$

From (12), we can see that the Xie-Beni index loses its ability to validate (U, V) pairs from the FCM for the large value of c .

The proposed index, $c \rightarrow n$: Since (11) holds for this case, we have

$$\lim_{c \rightarrow n} v_{XB}(U, V; X) = \lim_{c \rightarrow n} \frac{\sum_{j=1}^n \sum_{i=1}^c u_{ij}^2 \|x_j - v_i\|^2 + \frac{1}{c} \sum_{i=1}^c \|v_i - \bar{v}\|^2}{n \left[\min_{i \neq k} (\|v_i - v_k\|^2) \right]}$$

$$= \frac{\frac{1}{n} C_X}{\min_{i \neq k} (\|v_i - v_k\|^2)} \quad (13)$$

where C_X is the total scatter matrix of X . From (13), we can see that the proposed index keeps its ability to validate (U, V) pairs from the FCM for large value of c . Here, we do not discuss on the intuitive meaning and mathematical justification of the proposed index, which are required for the any new validity functional, because the research on those is on the way.

3. Numerical examples on cluster validity

In this section, we consider three examples of data sets to show the effectiveness of the proposed cluster validity. We first present a simple example with $c = 2$ as the preferred clusters, which is known as the butterfly data set [1] to provide insights into the limiting behavior of the cluster validity indexes. We then present two examples which are the derivatives of the butterfly data set and have $c = 3$ and $c = 4$ as the preferred clusters, respectively.

Example 1: We consider the butterfly data set X_1 of 15 data points in $p = 2$ dimensions shown in Fig.1, and listed in Table 1. Data points (2,2), (3,2), and (4,2) form a bridge or neck between the wings of the butterfly. Another interpretation of the pattern is that points in the wings were drawn from two fairly distinct classes; points in the neck are noise.

Example 2: We consider a set X_2 of 22 data points in

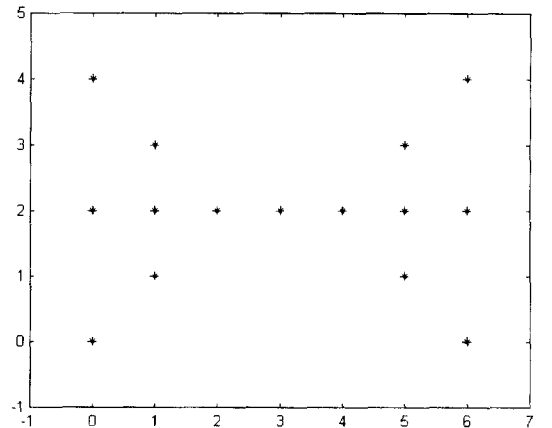


Fig. 1. Example 1: X_1

Table 1. Data sets X_1, X_2, X_3

Points	Data set X_1 ($c=2$)		Data set X_2 ($c=3$)		Data set X_3 ($c=4$)	
	x	y	x	y	x	Y
1	0	0	0	0	0	0
2	0	2	0	2	0	2
3	0	4	0	4	0	4
4	1	1	1	1	1	1
5	1	2	1	2	1	2
6	1	3	1	3	1	3
7	2	2	2	2	2	2
8	3	2	3	2	3	2
9	4	2	4	2	4	2
10	5	1	5	1	5	1
11	5	2	5	2	5	2
12	5	3	5	3	5	3
13	6	0	6	0	6	0
14	6	2	6	2	6	2
15	6	4	6	4	6	4
16			1	6	1	6
17			2	5	2	5
18			3	4	3	4
19			3	5	3	5
20			3	6	3	6
21			4	5	4	5
22			5	6	5	6
23					1	-2
24					2	-1
25					3	0
26					-1	27

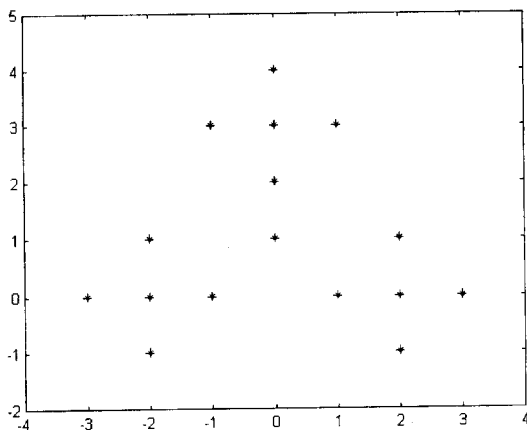


Fig. 2. Example 2: X_2

$p = 2$ dimensions shown in Fig. 2, and listed in Table 1. A data point (0, 1) forms a bridge among three diamonds. Another interpretation of the pattern is that points in the each diamond were drawn from three fairly distinct classes; a point s ; a point in the neck is noise.

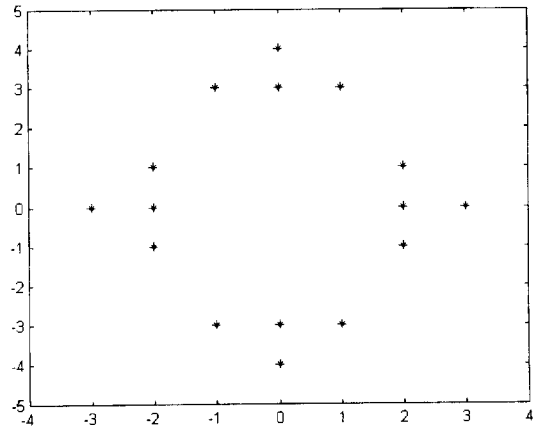
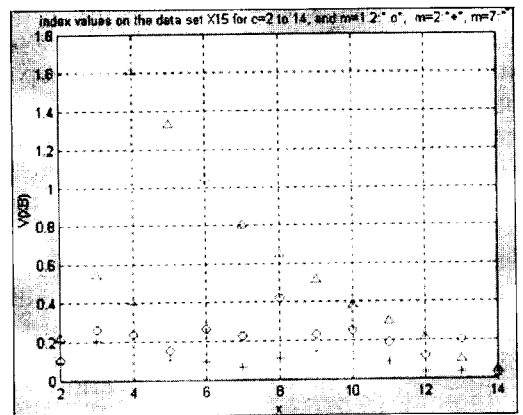
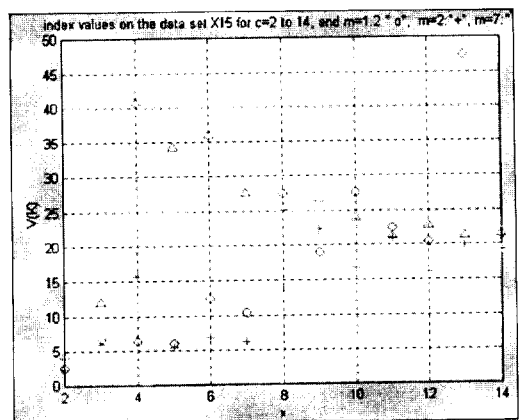


Fig. 3. Example 3: X_3



(a) The Xie and Beni index v_{XB}



(b) The proposed index v_K

Fig. 4. Index values on the data set X_1 for $c=2$ to 14, and $m=1.2, 2, 7$

Example 3: We consider a the butterfly data set X_3 of 29 data points in $p = 2$ dimensions shown in Fig. 3, and

Table 2. Values of c chosen by each index for the data sets X_2 and X_3

m	$X_2 : c^*=3$		$X_3 : c^*=4$	
	v_{XB}	v_K	v_{XB}	v_K
1.2	15	3	15	4
2.0	15	3	15	4
3.0	15	3	15	4
4.0	15	3	15	4
5.0	15	3	15	4
6.0	15	3	15	4
7.0	15	3	15	4

listed in Table 1. Data points in each triangular were drawn from four fairly distinct classes.

For each of data sets shown in the above we performed the FCM with the terminating criterion $\epsilon = 1.0e-8 \geq \|U_i - U_{i-1}\|$ for different weighting exponents $m = 1.2, 2, 3, 4, 5, 6$ and 7 , and $c = 2, 3, \dots, n-1$. Fig. 4 shows index values by each of the Xie-Beni index and the proposed index on the data set X_1 for $c = 2$ to 14 , and $m = 1.2, 2$, and 7 .

The preferred value of c for the tested data is 2 . From Fig. 4, we can see that the proposed index correctly points to the preferred value of c for each weighting exponent, but the Xie-Beni index points to $c = 14$. This behavior is consistent with the fact, which is the Xie-Beni index loses its ability to validate (U, V) pairs from the FCM for the large value of c , discussed in the previous section. Table 2 lists the value of the number of clusters chosen by each of the Xie-Beni index and the proposed index.

Since the preferred values of c are 3 and 4 , respectively, we see that the proposed index correctly points to the preferred values $c = 3$ and $c = 4$ for each weighting exponent but the Xie-Beni index points to $c = 15$ in every case. From these results, we conclude that the proposed cluster validity index shows the superior performance to the Xie-Beni index, and the Xie-Beni index may be unreliable.

4. Conclusions

In this paper, we have proposed a cluster validation index to eliminate the monotonically decreasing

tendency, which is the typical flaw of the conventional cluster validity indexes, when the number of cluster gets very large and close to the number of data points. We have reviewed the FCM algorithm and some cluster validity criteria, and discussed on the limiting behavior of the proposed validity index. Finally, numerical examples showing the effectiveness of the proposed cluster validity index have been provided.

Researches on the description of intuitive meaning, the mathematical justification, and applications to real data sets are on the way.

References

- [1] J. C. Bezdek, *Pattern Recognition with Fuzzy Objective Function Algorithms*, Plenum, New York, 1981.
- [2] G. W. Milligan, "Clustering validation," in *Clustering and Classification*, P. Arabie, L. J. Hubert and G. D. Soete, Ed. World Scientific, Singapore, 1996.
- [3] R. Krishnappuram and J. M. Keller, "A possibilistic approach to clustering," *IEEE Trans. Fuzzy Syst.*, Vol. 1, No. 2, pp. 98-110, 1993.
- [4] N. R. Pal, K. Pal and J. C. Bezdek, "A mixed c-means clustering model," in *Proc. FUZZ-IEEE'97*, pp. 11-21, 1997.
- [5] J. C. Dunn, "Indices of partition fuzziness and the detection of clusters in large data sets," in *Fuzzy Automata and Decision Processes*, M. M. Gupta, Ed. Elsevier, New York, 1976.
- [6] Y. Fukuyama and M. Sugeno, "A new method of choosing the number of clusters for the fuzzy c-means method," in *Proc. 5th Fuzzy Syst. Symp.*, pp. 247-250, 1989 (in Japanese).
- [7] J. C. Bezdek and N. K. Pal, "Some new indexes of cluster validity," *IEEE Trans. Systems, Man, and Cyber.-Part B*, Vol. 28, No. 3, pp. 301-315, 1998.
- [8] X. L. Xie and G. A. Beni, "Validity measure for fuzzy clustering," *IEEE Trans. Pattern and Machine Intell.*, Vol. 3, No. 8, pp. 841-846, 1991.
- [9] N. K. Pal and J. C. Bezdek, "On cluster validity for the fuzzy c-means model," *IEEE Trans. Fuzzy Syst.*, Vol. 3, No. 3, pp. 370-379, 1995.
- [10] M. Sugeno and T. Yasukawa, "A fuzzy logic based approach to qualitative modeling," *IEEE Trans. Fuzzy Syst.*, Vol. 1, No. 1, pp. 7-31, 1993.
- [11] S. H. Kwon, "Cluster validity index for fuzzy clustering," *Electronics Letters*, Vol. 34, No. 22, pp. 2176-2177, 1998.



권 순 학 (Soon H. Kwon)

1983년 : 서울대학교 제어계측공학과
졸업(공학사)
1985년 : 서울대학교 대학원 제어계측
공학과 졸업(공학석사)
1984년~1986년 : 삼성전자(주)주임연구원
1986년~1991년 : 한국과학연구원 연구원
1995년 : 동경공업대학 졸업(공학박사)

1996년~현재 : 영남대학교 전기전자공학부 조교수
관심분야 : 뉴시, 지능시스템의 모델링 및 제어



이 해 영 (Haiyoung Lee)

1984년 : 부산대학교 전기기계공학과
졸업(공학사)
1986년 : 한국과학기술원 전기 및 전자
공학과 졸업(공학석사)
1990년 : 한국과학기술원 전기 및 전자
공학과 졸업(공학박사)
1990년~1993년 : 포항종합제철(주) 기술
연구소 연구원

1993년~현재 : 영남대학교 전기전자공학부 조교수
주관심분야: 산업공정자동화



최 익 (Ick Choy)

1979년 : 서울대학교 전기공학과 졸업
(공학사)
1981년 : 서울대학교 대학원 전기공학과
졸업(공학석사)
1990년 : 서울대학교 대학원 전기공학과
졸업(공학박사)
1982년~현재 : 한국과학기술연구원 지능
제어연구센터장, 책임연구원

주관심분야: 전력전자, 지능제어