

# A New Spatial Fuzzy C-Means for Spatial Clustering

Yingdi Guo, Kunhong Liu\*, Qingqiang Wu, Qingqi Hong, Haiying Zhang

Software Department of Software School

Xiamen University

Xiamen, Fujian Province, China

578078128@qq.com; lkhqz@163.com; wuqq@xmu.edu.cn; hongqq@xmu.edu.cn; zhang2002@xmu.edu.cn

*Abstract:* Fuzzy C-means is a widely used clustering algorithm in data mining. Since traditional fuzzy C-means algorithms do not take spatial information into consideration, they often can't effectively explore geographical data information. So in this paper, we design a Spatial Distance Weighted Fuzzy C-Means algorithm, named as SDWFCM, to deal with this problem. This algorithm can fully use spatial features to assign samples to different clusters, and it only needs to calculate the memberships one time, which reduces the running time greatly compared with other spatial fuzzy C-means algorithms. In addition, we also propose two new criteria, named as DESC and PESC, for evaluating spatial clustering results by measuring spatial and regular information separately. The experiments are carried out based on real petroleum geology data and artificial data, and the results show that SDWFCM can achieve better performance compared with traditional clustering method, and our spatial cluster indices can provide the assessment of clusters by taking spatial structure into consideration effectively.

*Key-Words:* spatial clustering, fuzzy c-means, evaluation criteria.

## 1. Introduction

Spatial clustering is an important research field of data mining, and it has been widely used in geography, geology, remote sensing, mapping and other disciplines. When clustering spatial data, each sample is divided into two parts: the spatial information and general properties. But in most cases, in the task of clustering spatial data, researchers mainly consider general attributes, with leaving the spatial information somehow ignored. However, for geological data, it is obvious that simply considering regular attributes can not effectively reflect the characteristics of sample data. So some studies have been done to combine the spatial and non-spatial attributes. For example, X.Y. Li et al. proposed the space coordinate integration

[1], and G.Q. Li et al. proposed a spatial clustering algorithm based on dual distance [2].

Traditional fuzzy C-means (FCM) algorithm has been used in different research fields. It allocates samples' membership according their probability of belonging to a cluster. Because of its effectiveness, more and more methods are proposed to improve it. Fan et al [3] introduced Suppressed FCM algorithm. It increases speed of FCM by prizing the biggest membership value and suppressing the others. However, when membership values are close, simply rewarding the biggest membership value may not be reasonable. F. Zhao et al. [4] proposed optimal-selection based suppressed fuzzy c-means clustering algorithm with self-tuning non local spatial information. It constructs

gray level histogram firstly, and then an optimal-selection based suppressed strategy is applied to improve FCM, which makes a successful application on image segmentation. Zexuan Ji et al. added weighted image patch in FCM algorithm [5]. It replaces pixels with image patches as the basic unit of clustering, which can reduce the impact of noise and saving times. It works well on brain MR images. Krishna Kant Singh et al. [6] presented a neuro fuzzy clustering algorithm for classification of remote sensing images. Hongbao Cao et al. [7] proposed IAFCM algorithm used a new objective function with a different regulation term. And it appears to be more effective in controlling the shape of the gain field. It has been successfully applied to the classification of M-FISH images. M-FISH image is a combinatorial labeling technique that is developed for the analysis of human chromosomes. S. Krinidis and V. Chatzis [8] proposed a robust FCM framework for image clustering, named as FLICM. It changes the conventional FCM objective function by adding local spatial and gray level information. M.G. Gong et al. [9] proposed KWFLICM based on FLICM, which adds a trade-off weighted fuzzy factor to the objective function and kernel method. The adaptive trade-off weighted fuzzy factor depends on the local spatial constraint and local gray-level constraint. Hesam Izakian et al. [10] revisit and augment the FCM to make it applicable to spatiotemporal data by discussing an augmented distance function. Some researchers also applied FCM to combine with evolutionary based algorithms to real world problems [11]. The survey of FCM based methods are given in [12].

However, traditional FCM algorithms are not proper for dealing with spatial data because such algorithms can't separate spatial features with general features. Based on this consideration, in recent years, many researchers have put forward some schemes to solve this problem. For example, C.P. Hu et al. [13] and Chuang et al [14] proposed two similar Spatial Fuzzy C-means algorithms,

which named as SFCM and sFCMpq, respectively. These algorithms use two different spatial functions for smoothing membership value, and use parameters to balance the importance of the normal membership and the spatial function's value. Y.Y. Wang et al. [15] used some methods to further improve sFCMpq. However, the above methods need calculate the membership degree twice, which is quite time-consuming when dealing with large-scale problems. And the process of smoothing membership degree would tend to unavoidably blur the characters of samples.

This paper presents a new Spatial Distance Weighted Fuzzy C-Means algorithm, named as SDWFCM. It can directly calculate weighted distance, and the memberships are only calculated once. So it is much faster than other methods in application. In addition, to evaluate the effectiveness of algorithms, we also propose two new measurements for spatial clustering. The experimental results show that SDWFCM can be more efficient and effective than other methods, and the evaluation criterion is qualified to be new indicators for the comparisons of spatial clustering algorithms.

## 2. Method

### 2.1 Spatial distance weighted based fuzzy C-means

Traditional FCM was proposed by Bezdek J.C [13]. It divides data set  $X = \{x_1, x_2, \dots, x_n\}$  into  $c$  fuzzy clusters by minimizing the cost function, defined as formula (1):

$$J(X; U, V) = \sum_{i=1}^c \sum_{k=1}^n (u_{ik})^m \|x_k - v_i\|^2 \quad (1)$$

Where  $v_i$  represents the center of cluster  $i$ ,  $x_i \in R^d$ .  
 $U = [u_{ik}]$  is a membership matrix of dataset, and it satisfied formula (2):

$$\sum_{i=1}^c u_{ik} = 1, \quad 1 \leq k \leq n, \quad u_{ik} \in [0,1] \quad (2)$$

Here,  $m \in [1, \infty]$ , and  $m$  is a weighted index used to defines the degree of fuzzy classification. When  $m=1$ , the method degradation to the hard clustering. And usually we set  $m=2$ .

$$u_{ik} = \frac{1}{\sum_{j=1}^c (d_{ik} / d_{jk})^{1/(m-1)}}, \quad 1 \leq i \leq c, 1 \leq k \leq n \quad (3)$$

$$v_i = \frac{\sum_{k=1}^n (u_{ik})^m x_k}{\sum_{k=1}^n (u_{ik})^m}, \quad 1 \leq i \leq c \quad (4)$$

FCM algorithm iteratively computes the cluster centers and the corresponding memberships according to formula (3) and (4), until the clustering centers vary within a pre-defined range or the maximum number of iteration reaches.

## 2.2 SFCM

Although FCM works well in many real-world applications, it can't be applied to deal with spatial data directly, because there are two types of features in spatial clustering problems. The first type of features is regular feature, which refers to ordinary property of non-spatial data. Such typical features are age, income, and et al. The other type of features is directly referenced to a location of the earth, such as Latitude and Longitude. They can be represented by coordinates or vectors. In many applications, researchers may ignore the specialty of spatial features, and just use traditional clustering algorithms by regarding spatial features as regular

features. However, it may unavoidably lead to the loss of some characters of samples and the spatial structures among samples. To solve this problem, some researchers devote to the design of spatial clustering algorithms. For example, C.P. Hu et al [9] proposed SFCM by defining a spatial function as follow:

$$h_{ij} = \sum_{s_k \in NB(S_j)} u_{ik} \quad (5)$$

where  $NB(S_j)$  represents the neighborhood object set of a spatial object  $S_j$ . It is instinctive that when most of  $S_j$ 's neighborhood objects belong to the same cluster, the probability of  $S_j$  belonging to this cluster would be high. Therefore, the algorithm made the following modifications to the membership matrix:

$$u'_{ij} = \frac{u_{ij}^p h_{ij}^q}{\sum_{i=1}^c u_{kj}^p h_{kj}^q} \quad (6)$$

where  $p$  and  $q$  act as weighted parameters, and both are used to control the weight of  $u_{ij}$  and  $h_{ij}$ . When  $p=1$  and  $q=0$ , SFCM degrades into the traditional fuzzy C-means algorithm.

Although SFCM can obtain more reasonable results compared with traditional FCM, there are still some shortcomings in this algorithm. Firstly, the value of  $p$  and  $q$  can't define directly, and no relationship can be found between  $p$  and  $q$ . So in practical applications, it's hard to know how to set  $p$  and  $q$ . And the worse of all, there is no way to find out whether these two parameters should be increased or decreased when tuning the parameters of the algorithm. Secondly, because the memberships need to be calculated twice, the speed of SFCM is

very slow in case of large scale datasets. So our algorithm aims to tackle these problems.

### 2.3 SDWFCM

A Spatial Distance Weighted based FCM (SDWFCM) is proposed to solve the problems mentioned above. The principle of this algorithm is described as follows.

First of all, SDWFCM and traditional FCM are consistent in computing the center of cluster, as shown in the following formula:

$$v_i^{(t)} = \frac{\sum_{k=1}^n (u_{ik}^{(t-1)})^m x_k}{\sum_{k=1}^n (u_{ik}^{(t-1)})^m}, 1 \leq i \leq c \quad (7)$$

Where the  $v_i^{(t)}$  represents the  $i$ -th cluster center after the  $t$ -th iteration,  $u_{ik}$  indicates the membership of the  $k$ -th sample belonging to the  $i$ -th cluster.

Before calculating the distance between sample and the cluster center, we define a spatial distance weighted function as follows:

$$f_{ij} = \frac{\sum_{s_k \in NB(s_j)} d_{ik}}{\text{Min}\{\sum_{s_k \in NB(s_j)} d_{lk}, 1 \leq l \leq c\}} \quad (8)$$

This function is employed to modify distances as follows:

$$D_{ij} = (1 - \lambda) \cdot d_{ij} \cdot f_{ij} + \lambda \cdot d_{ij} \quad (9)$$

where  $\lambda$  is the weighted coefficient, which ranges within (0, 1). And  $d_{ij}$  is the original Euclidean distance of the sample to the cluster center.  $D_{ij}$  is a weighted distance that sample to the cluster center, and it is deployed in the process of reassigning memberships.

$$u_{ik}^{(t)} = \frac{1}{\sum_{j=1}^c (D_{ik} / D_{jk})^{1/(m-1)}} \quad (10)$$

In this way, the weight of spatial features is fully controlled by  $\lambda$ . When  $\lambda$  is set close to 1, the spatial information will play an important role in the final result. At the same time, the membership only needs to be calculated once, leading to save much more time compared with SFCM. The algorithm flow is described as below.

Let  $U$  be the membership matrix, and its elements represent the membership of  $k$ -th sample belonging to  $i$ -th cluster. Matrix  $D$  records the Euclidean distance of different samples to cluster centers, where  $d_{ij}$  represents Euclidean distance of the  $j$ -th sample to the  $i$ -th cluster center. They are calculated using ordinary feature, just like the standard clustering algorithms.  $V$  is the matrix of clustering center, and the elements of  $v_i^t$  represent the  $i$ -th cluster's center after the  $t$ -th iteration. The work flow of SDWFCM is described as below.

1. Initializing membership matrix  $U$ ;
2. Calculating the cluster center matrix  $V$  according to formula (7);
3. Calculating the Euclidean distance matrix  $D$  based on regular features;
4. Calculating weighted distance  $D_{ij}$  according to formula (8) and (9);
5. Calculating the new membership degree based on formula (10);
6. If the number of iterations has reached the maximum run time, or the variance of cluster centers meets a criteria, the program stops; otherwise, it repeats step 2-5.

### 3. Evaluation index

#### 3.1 Current situation of clustering index

External clustering evaluation is an effective way for the validity of clustering results, and effective criterion is key to fairly comparing the performances of different clustering algorithms or tuning parameters a clustering algorithm.

Traditional clustering evaluation measurements mainly concern about distance between classes or within a class. For example, Dunn's Index to calculate the ratio of minimum distance of samples in a same cluster and the maximum distance between clusters. Davies-Bouldin's Index compute the average distance within clusters and the maximum distance between clusters. CS Index evaluates the cluster

results in a similar way by considering distance [17]. But these indices only calculate the distance in regular feature sets. When evaluating spatial clustering results, such as geographic data, these evaluation measures are insufficient due to lacking of measures for spatial structure. Chunchun Hu et al. [18] proposed the IFV evaluation index for spatial fuzzy clustering, but the evaluation algorithm is of high computation complexity, and can't reflect the spatial structure of geographic data. Q.L. Liu et al. [19] put forward a spatial clustering evaluation measurement based on field theory, but it is also time consuming.

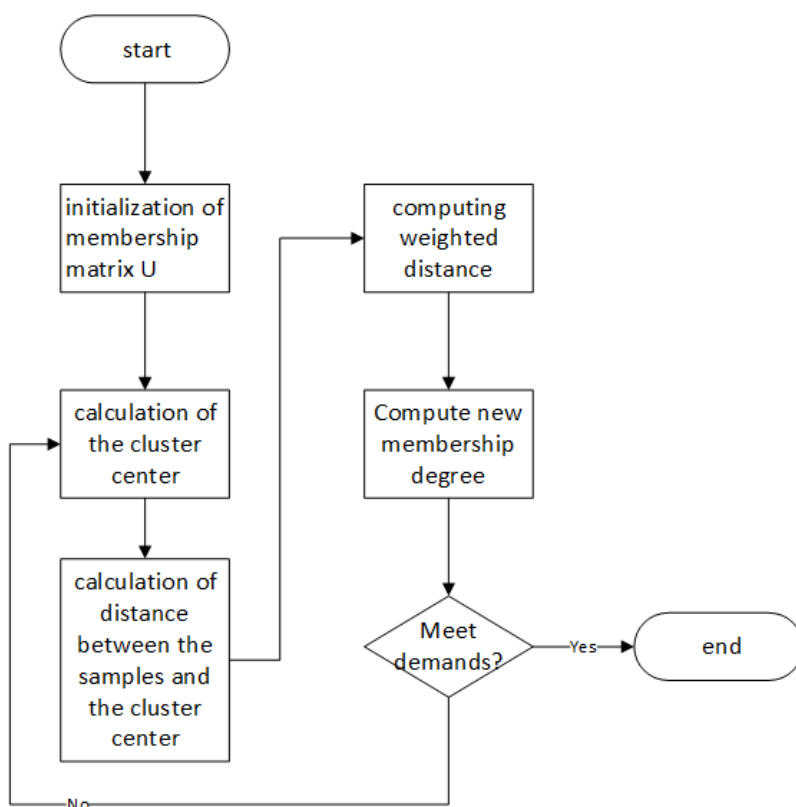


Fig. 1. Flowchart of DWSFCM algorithm

### 3.2 Two new clustering evaluation indices

Two new evaluation criterions for clustering spatial data are proposed here: discrete evaluation for spatial clustering (DESC) and proximity evaluation for spatial clustering (PESC).

#### 1. DESC

This criterion is designed by this assumption: a clustering algorithm should assign as much as possible geographically adjacent samples to a same cluster in the case that they are also neighbors in regular feature subspace.

Let  $Block_i$  denote the  $i$ -th sets in which the geographically adjacent samples are assigned to the same cluster. It is called as the  $i$ -th spatially connected block. The samples in the set must be in a same cluster, and there is at least one sample in a set. The connected blocks form a bigger set, named as  $Block$ . Assume that the number of all the connected blocks is  $N_b$  for  $Cluster_k$  in a set, containing all connected blocks belonging to cluster  $k$ . Let  $Cn_k$  represents the number of samples belonging to  $k$ -th cluster. Let  $v_i$  denote the number of samples in  $Block_i$ , or the size of the  $i$ -th collection. And  $V_i$  is the size of the  $i$ -th collection divided by the number of samples belonging to the  $k$ -th cluster, which can be described as:

$$V_i = \frac{v_i}{Cn_k} \quad \{Block_i \in Cluster_k\} \quad (11)$$

Let  $\bar{d}_i$  represent the average distance of all samples in the collection to the collection center for general attributes. The formula of Discrete Evaluation for Spatial Clustering can be described as follows:

$$f_{desc} = \sum_{i=1}^{N_b} \frac{V_i^2}{d_i^2} \quad (12)$$

where  $N$  represents the total number of all samples in the sample. When the  $v_i$ 's value is 1, we let the  $\bar{d}_i$  be 1. From the formula, it can be found that the bigger size a connecting block is, the higher score it obtains. At the same time, if the average distance of samples in a connected block to the block center is greater, the score is lower. In addition, when the number of connected block is 1, the average distance of  $\bar{d}_i$  is 0, which indicates an illegal situation. In such case,  $\bar{d}_i$  is simply set to 1.

#### 2. PESC

This index show how close the samples in the same cluster are. It is defined as follows:

Let  $Block_i$  denote the  $i$ -th sets in which the adjacent samples are belong to a same cluster. It is named as the  $i$ -th connected block. The samples in the set must be in the same cluster, and each sample has at least one adjacent sample within the set, or the set contains only one sample.  $Cluster_k$  is a set, and it contains all connected blocks those belong to cluster  $k$ . Let  $Cn_k$  represent the number of samples which belong to  $k$ -th cluster, and  $Bn_k$  denote the number of blocks belonging to the  $k$ -th cluster. Let  $v_i$  denote the number of samples in  $Block_i$  or the size of the collection,  $V_i$  is the size of the  $i$ -th collection divided by the number of samples belonging to  $k$ -th cluster, as shown in (13).

$$V_i = \frac{v_i}{Cn_k} \quad \{Block_i \in Cluster_k\} \quad (13)$$

The calculation formula of Proximity Evaluation for Spatial Clustering is:

$$f_{pesc} = \sum_{k=1}^c \sum_{D_{ij}} \frac{V_i \cdot V_j}{d_{ij} \cdot Bn_k^2} \quad Block_i, Block_j \in Cluster_k \quad (14)$$

In (14),  $D_{ij}$  represents the spatial distance of the closest sample pair belonging to  $Block_i$  and  $Block_j$ , respectively, and  $d_{ij}$  represents the distance of two block centers for regular attributes. So the closer of connected blocks of the same cluster are, the higher the scores they obtain. And the more samples in the same connected block, the higher score. And it should be noted that if there is a large block in  $Cluster_i$  with a great number of small blocks, it still gets a low score.

## 4. Experiments and analysis

### 4.1 Real oil seismic exploration data

We compare the performance of three different algorithms: K-means (with default parameter settings in Matlab toolbox), SFCM, and DWSFCM. The number of clusters is set to 7 in all experiments. The visualized results for three algorithms are shown as Fig. 2- Fig. 4.

Comparing the clustering results, we can see that the results are all filled with small segments, but the results produced by Kmeans are shown to be the most scattered. Too many trivial segments indicate that clustering algorithm fails to exploit the spatial information reasonably and sufficiently. The results produced by SFCM and DWSFCM are much better as shown in Fig. 2 and 3. In addition, it is necessary to compare the evaluation index achieved by these algorithms. And we use some evaluation indices, Iidx, CDVM, Dunn's Index, DESC proposed in this paper to evaluate the results of clustering under different parameter. We don't use PESC proposed in this paper because the scale of oil seismic exploration data is so huge that the task of finding the closest sample pair for blocks become a tough challenge. And we can't obtain results in a reasonable time, so we have to ignore the PESC index in real data set. The results are

shown in the Table 1. It should be noted that all results are unitless.

For all indices listed in Table I, the bigger, the better. From Table I, we can observe that SFCM algorithm beats K-means in Iidx scores, but its CDVM and Dunn scores are worse than K-means. For DESC, when the SFCM parameter  $p$  value is higher, the results are better than K-means. But the SFCM algorithm's parameters do not reflect what role the weigh of the spatial information plays in clustering.

In contrast, DWSFCM's parameter  $\lambda$  can directly reflect the importance of spatial information. When  $\lambda = 0.5$ , DWSFCM gets the best scores compared with other two methods in the first three indices. And when  $\lambda = 0.4$ , it gets the highest spatial index in DESC, and the value is much higher than those of K-means and SFCM with different parameter settings. The large score of DESC index shows that less noise and larger blocks in clustering results. When  $\lambda = 0.5$ , its DESC score is also better than other algorithms. Thus, the DWSFCM algorithm can achieve better performance compared to traditional K-means and SFCM.

The success of DWSFCM lies in that its parameters directly reflect the importance of spatial information. It can be found that, with the increase of  $\lambda$ , its results are getting better at first, and then get worse. So the proportion of spatial information can be neither blindly increased, nor ignored. When  $\lambda$  is close to 0.5, the results reach an optimal situation with treating the spatial information and attribute information equally. In such case, we can find that DWSFCM achieve the best performance.

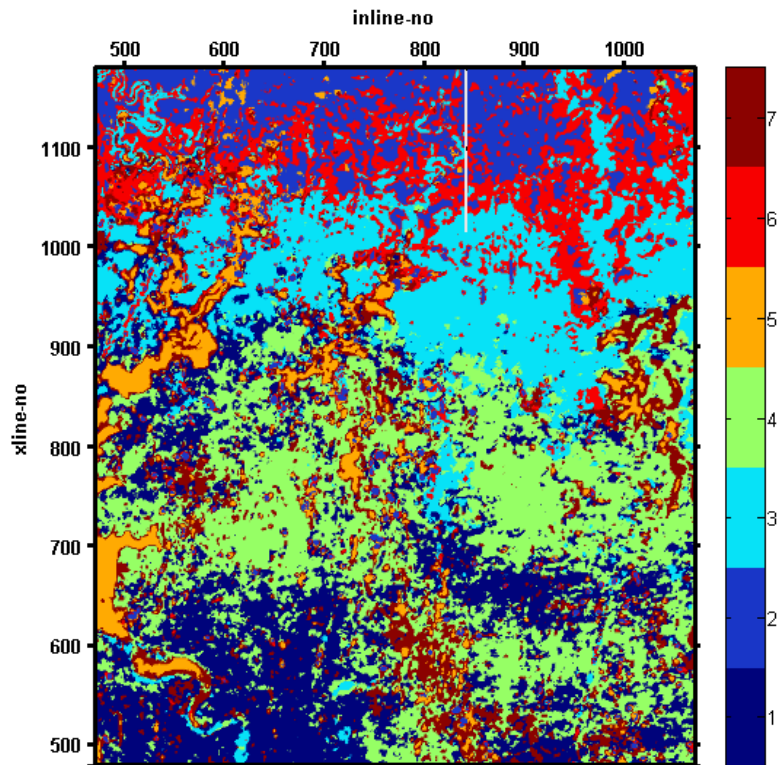


Fig. 2. K-Means results

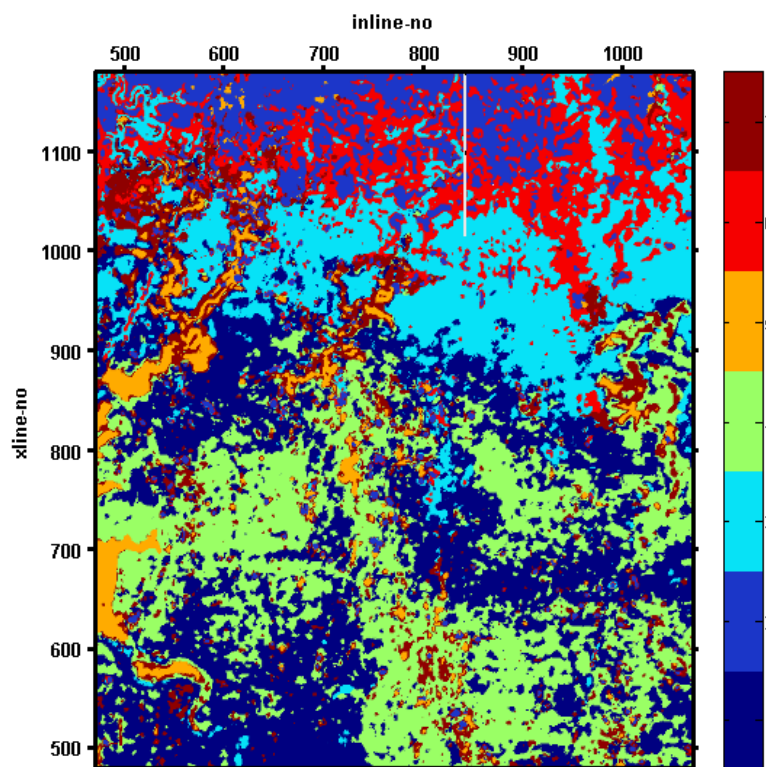


Fig. 3. SFCM results ( $p = 5, q = 2$ )



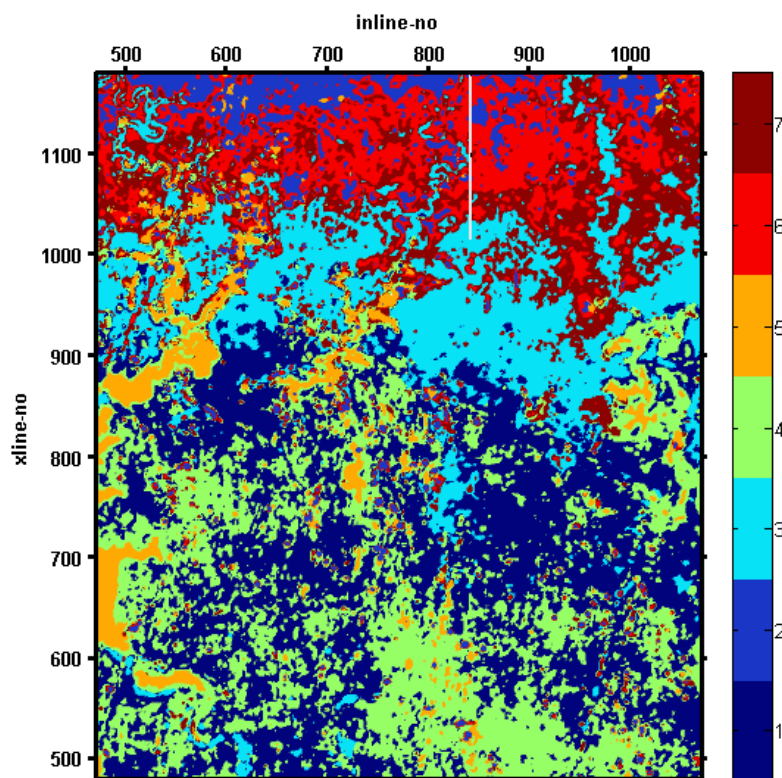
Fig. 4. DWSFCM results ( $\lambda = 0.5$ )

TABLE I Results of different algorithms based on real world spatial data

	<b>IIdx</b>	<b>CDVM</b>	<b>Dunn</b>	<b>DESC</b>
<b>K-Means</b>	0.0092252	0.5208984	0.0182718	53.9333421
<b>SFCM (p=5,q=2)</b>	0.0137251	0.5192598	0.0076759	47.3330820
<b>SFCM (p=8,q=2)</b>	0.0122435	0.5205313	0.0108863	53.0173389
<b>SFCM (p=10,q=2)</b>	0.0114494	0.5205647	0.0123926	53.2904940
<b>SFCM (p=12,q=2)</b>	0.0108670	0.5206423	0.0134363	55.8093210
<b>SFCM (p=26,q=2)</b>	0.0098918	0.5208213	0.0160930	55.0747204
<b>DWSFCM (<math>\lambda = 0.3</math>)</b>	0.0090212	0.5226890	0.0241802	70.4694010
<b>DWSFCM (<math>\lambda = 0.4</math>)</b>	0.0103273	0.5224728	0.0186201	74.4526495
<b>DWSFCM (<math>\lambda = 0.5</math>)</b>	0.0165571	0.5237769	0.0178778	55.5099774
<b>DWSFCM (<math>\lambda = 0.6</math>)</b>	0.0130327	0.5219415	0.0109377	44.1307562
<b>DWSFCM (<math>\lambda = 0.7</math>)</b>	0.0093480	0.5225258	0.0108132	59.3899967

## 4.2 Artificial data

In order to make the comparison of algorithms more clearly, we set up an artificial data. The artificial data has 10000 data points, these data points distributed evenly in the  $100 * 100$  dot matrix. And Fig. 5 illustrates the original data by using each pixel to represent a sample in

the corresponding spatial location. In the following experiments, the results are all illustrated in the same way. And it should be noted that these figures are unitless. In our experiments, each sample contains 4 attributes, produced according to the following rules.

- 1) if the data points corresponding to the pixel as red, attribute 1 would be a

random numbers 0.3 to 0.8, attribute 2 would be a random number from 0.6 to 0.9, attribute 3 would be a random numbers from 0.1 to 0.6, attribute 4 would be a random numbers 0.4 to 0.9.

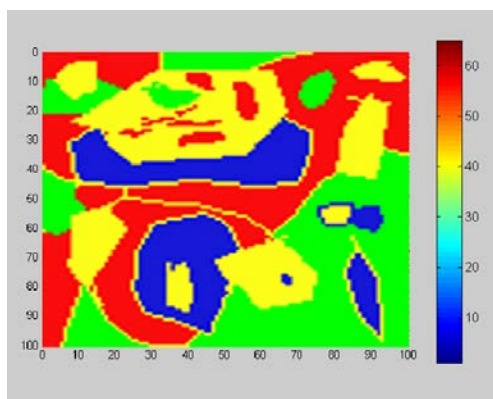


Fig. 5. The artificial data

- 2) if the data points corresponding to the pixel as yellow, attribute 1 would be a random numbers 0.0 to 0.7, attribute 2 would be a random number from 0.4 to 0.8, attribute 3 would be a random numbers from 0.2 to 0.5, attribute 4 would be a random numbers 0.6 to 1.0.
- 3) if the data points corresponding to the pixel as green, attribute 1 would be a random numbers 0.5 to 0.8, attribute 2 would be a random number from 0.1 to 0.8, attribute 3 would be a random numbers from 0.2 to 0.8, attribute 4 would be a random numbers 0.2 to 0.7.
- 4) if the data points corresponding to the pixel as blue, attribute 1 would be a random numbers 0.1 to 0.5, attribute 2 would be a random number from 0.2 to 0.6, attribute 3 would be a random numbers from 0.4 to 0.9, attribute 4 would be a random numbers 0.1 to 0.5.

In this way, an artificial data set is generated in size of  $10000 \times 4$ . The ideal clustering results should be able to produce the original artificial data source diagram. However, according to the

artificial data set generating rule, some samples in the data set will become outliers. And effective clustering algorithms should reduce the effect of outliers.

The corresponding clustering results of each algorithm are shown in Fig.6-Fig.8. From Fig. 6, we can see that although the K-means algorithm can correctly pick up samples in a same cluster and produce boundaries for different clusters clearly, it can not deal with noise, resulting with the clustering results spotted. And it is obvious that it can not get rid of the influence of random noise.

Fig.7 shows that SFCM algorithm can produce better and clean results compared with those of K-means in general. And there is less noise. And DWSFCM has the less noise among all results, as shown in Fig. 8.

We also use evaluation index to compare the algorithms with different parameters, as shown in Table II.

All larger values in Table II indicate better performance. From this table, we can see that the DWSFCM algorithm performs better than the other two algorithms judged by both DESC index and Dunn index. According to PESC index, DWSFCM algorithm beats K-means, but its PESC index is slightly inferior to that of SFCM algorithm, which may reveal that for the spatial data, parameter settings are very important.

On the other hand, because PESC is designed to measure the relationship between blocks in a same cluster, the results also indicate that the SFCM algorithm produce compact blocks in clusters when P is smaller. In all experiments, the performances of K-Means are relatively poor in considering both DESC index and PESC index. And Fig.6 also shows that its clustering results contain more noise.

TABLE I Results of different algorithms based on artificial spatial data

	<b>IIdx</b>	<b>CDVM</b>	<b>Dunn</b>	<b>DESC</b>	<b>PESC</b>
<b>K-Means</b>	0.04140	0.39425	0.15854	4023.72	236.81
<b>SFCM (p=3,q=2)</b>	0.05207	0.39195	0.14578	4062.48	431.01
<b>SFCM (p=5,q=2)</b>	0.05426	0.39548	0.14707	5477.70	315.15
<b>SFCM (p=7,q=2)</b>	0.05312	0.39608	0.14872	4948.03	328.03
<b>SFCM (p=11,q=2)</b>	0.05037	0.39598	0.15126	4894.24	322.98
<b>SFCM (p=23,q=2)</b>	0.04660	0.39551	0.15465	4205.57	248.14
<b>DWSFCM (<math>\lambda = 0.3</math>)</b>	0.03220	0.39133	0.16172	9055.64	382.86
<b>DWSFCM (<math>\lambda = 0.4</math>)</b>	0.03297	0.39059	0.16100	6986.30	293.14
<b>DWSFCM (<math>\lambda = 0.5</math>)</b>	0.03404	0.38999	0.16064	6026.37	278.86
<b>DWSFCM (<math>\lambda = 0.6</math>)</b>	0.03375	0.38883	0.15961	6044.38	284.35
<b>DWSFCM (<math>\lambda = 0.7</math>)</b>	0.03367	0.38777	0.15885	5234.39	259.03

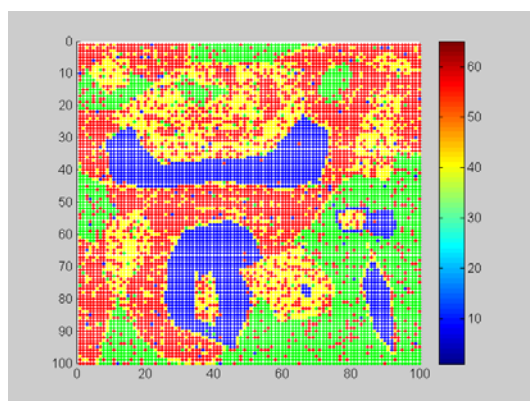
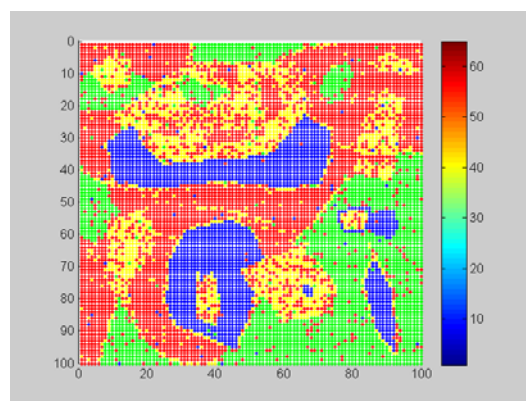
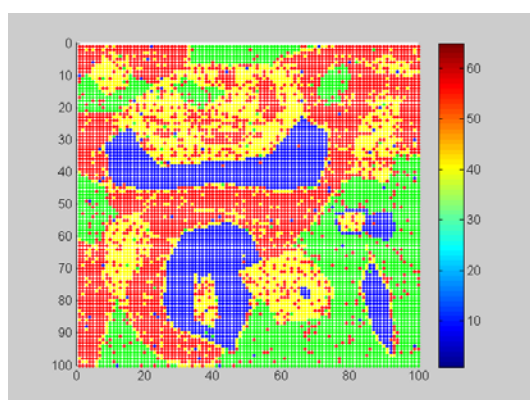


Fig. 6. K-Means results

Fig. 8. DWSFCM results ( $\lambda = 0.5$ )Fig. 7. SFCM results ( $p = 5, q = 2$ )

Although DWSFCM algorithm is worse than the other two algorithms when taking IIDx and CDVM indices into consideration, by comparing results in Fig. 6-Fig. 8, it can still be found that the results of DWSFCM are better than other two algorithms. In our opinion, this observation shows that the traditional indices are not appropriate for evaluating spatial cluster algorithms due to the ignoring of spatial structure. On the contrary, clustering evaluation index proposed in this paper can reflect the character of spatial information more clearly. As a conclusion, spatial information should be taken into account in the design of measures for spatial data mining algorithms.

## 5. Conclusions

In this paper, we propose a new fuzzy K-mean algorithm for clustering spatial data, named as DWSFCM. It can be applied to effectively deal with spatial features and regular features at the same time. To evaluate the results for spatial data, we also design two measurements, named as DESC and PESC. In our experiments, it is found that our method can usually achieve the best results compared with K-means and SFCM, and the indices we proposed can also help to give a new deep insight to spatial clustering results.

## Acknowledgements

This work is supported by National Science Foundation of China (No. 61100106 and 61303080), National Key S&T Project of China (No.2011ZX05004-003), and the Research Program of RIPED (No. 101002kt0b52135).

## References:

- [1] X. Li, X. Zheng, and H. Yan, On spatial clustering of combination of coordinate and attribute, *Geography and Geo-Information Science*, Vol. 20, No. 2, 2004, pp. 38–40.
- [2] G. Li, M. Deng, J. Zhu, and T. Cheng, A dual distance based spatial clustering method, *Acta Geodaetica et Cartographica Sinica*, Vol. 37, No. 4, 2008, pp. 482–488.
- [3] Zhao, Feng, Jiulun Fan, and Hanqiang Liu. Optimal-selection-based suppressed fuzzy c-means clustering algorithm with self-tuning non local spatial information for image segmentation, *Expert Systems with Applications*, Vol. 41, No. 9, 2014, pp. 4083-4093.
- [4] Ji, Zexuan, et al, Fuzzy c-means clustering with weighted image patch for image segmentation, *Applied Soft Computing*, Vol. 12, No. 6, 2012, pp. 1659-1667.
- [5] Singh, Krishna Kant, et al. A Fuzzy Kohonen Local Information C-Means Clustering for Remote Sensing Imagery, *IETE Technical Review*, Vol. 31, No. 1, 2014, pp. 75-81.
- [6] Cao, Hongbao, Hong-Wen Deng, and Yu-Ping Wang, Segmentation of M-FISH images for improved classification of chromosomes with an adaptive Fuzzy C-Means Clustering Algorithm, *IEEE TRANSACTIONS ON Fuzzy Systems*, Vol. 20, No. 1, 2012, pp. 1-8.
- [7] Stelios Krinidis and Vassilios Chatzis. A Robust Fuzzy Local Information C-Means Clustering Algorithm, *IEEE TRANSACTIONS ON IMAGE PROCESSING*, Vol. 19, No. 5, 2010, pp. 1328-1337.
- [8] Gong, Maoguo, et al., Fuzzy c-means clustering with local information and kernel metric for image segmentation, *IEEE Transactions on Image Processing*, Vol. 22, No. 2, 2013, pp. 573-584.
- [9] Izakian, Hesam, Witold Pedrycz, and Iqbal Jamal, Clustering Spatiotemporal Data: An Augmented Fuzzy C-Means, *IEEE Transactions on Fuzzy Systems*, Vol. 21, No. 5, 2013, pp. 855-868.
- [10] S. Mungle, L. Benyoucef, Y.J. SoN, M. K. Tiwari, A fuzzy clustering-based genetic algorithm approach for time-cost-quality trade-off problems: A case study of highway construction project, *Engineering Applications of Artificial Intelligence archive*, Vol. 26, No. 8, 2013, pp. 1953-1966.
- [11] G. Petersa, F. Crespoc, P. Lingrasd, R. Webere, Soft clustering – Fuzzy and rough approaches and their extensions and derivatives, *International Journal of Approximate Reasoning*, Vol. 54, No.2, 2013, pp 307–322.
- [12] H. Caiping, Research of key techniques on spatial data mining based on spatial autocorrelation, Ph.D. dissertation, Nanjing University of Aeronautics and Astronautics, 2007.

- [14] K.-S. Chuang, H.-L. Tzeng, S. Chen, J. Wu, and T.-J. Chen, Fuzzy c-means clustering with spatial information for image segmentation, *Computerized Medical Imaging and Graphics*, Vol. 30, No. 1, 2006, pp. 9-15.
- [15] W. Yuanyuan, The study of fuzzy c-means algorithm incorporating spatial information for brain MR image segmentation, Master's thesis, XIDIAN UNIVERSITY, 2012.
- [16] Bezdek James C., Robert Ehrlich, and William Full, FCM: The fuzzy c-means clustering algorithm, *Computers & Geosciences*, Vol. 10, No. 2, 1984, pp. 191-203.
- [17] L. Ying, The research on the method to measure the validity and to abstract knowledge of clustering, Master's thesis, Nanjing University of Aeronautics and Astronautics, 2008.
- [18] H. Chunchun, M. Lingkui, X. Wenjun, and Z. Xinzhong, Validity measure on fuzzy clustering for spatial data, *Geomatics and Information Science of Wuhan Univers*, Vol. 32, No. 8, 2007, pp. 740-743.
- [19] M. DENG, Q. LIU, G. LI, and T. CHENG, Field-theory based spatial clustering method, *Journal of Remote Sensing*, Vol. 14, No. 4, 2009, pp. 694-709.