

# A layered parallel algorithm for data clustering via fuzzy c-means technique

Yashuang Mu (✉ [muyashuang324@126.com](mailto:muyashuang324@126.com))

Henan University of Technology <https://orcid.org/0000-0003-0362-9640>

Wei Wei

Henan University of Technology

Hongyue Guo

Dalian Maritime University

Lijun Sun

Henan University of Technology

---

## Research Article

**Keywords:** Clustering, Fuzzy c-means, Map-Reduce, Parallel Computing

**Posted Date:** April 19th, 2021

**DOI:** <https://doi.org/10.21203/rs.3.rs-268783/v1>

**License:** © ⓘ This work is licensed under a Creative Commons Attribution 4.0 International License.

[Read Full License](#)

---

# A layered parallel algorithm for data clustering via fuzzy c-means technique

Yashuang Mu<sup>a,b,\*</sup>, Wei Wei<sup>c</sup>, Hongyue Guo<sup>d,e</sup>, Lijun Sun<sup>a</sup>

<sup>a</sup>Key Laboratory of Grain Information Processing and Control (Henan University of Technology),  
Ministry of Education, Zhengzhou 450001, P.R. China.

<sup>b</sup>School of Artificial Intelligence and Big Data, Henan University of Technology, Zhengzhou 450001,  
P.R. China.

<sup>c</sup>College of Information Science and Engineering, Henan University of Technology, Zhengzhou 450001,  
P.R. China.

<sup>d</sup>School of Maritime Economics and Management, Dalian Maritime University, Dalian 116026, China

<sup>e</sup>Collaborative Innovation Center for Transport Studies, Dalian Maritime University, Dalian 116026,  
China

---

## Abstract

In this study, a layered parallel algorithm via fuzzy c-means (FCM) technique, called LP-FCM, is proposed in the framework of Map-Reduce for data clustering problems. The LP-FCM mainly contains three layers. The first layer follows a parallel data partitioning method which is developed to randomly divide the original dataset into several subdatasets. The second layer uses a parallel cluster centers searching method based on Map-Reduce, where the classic FCM algorithm is applied to search the cluster centers for each subdataset in Map phases, and all the centers are gathered to the Reduce phase for the confirmation of the final cluster centers through FCM technique again. The third layer implements a parallel data clustering method based on the final cluster centers. After comparing with some famous classic random initialization sequential clustering algorithms which include K-means, K-medoids, FCM and MinMaxK-means on 20 benchmark datasets, the feasibility in terms of clustering accuracy is evaluated. Furthermore, the clustering time and the parallel performance are also tested on some generated large-scale datasets for the parallelism.

*Keywords:* Clustering, Fuzzy c-means, Map-Reduce, Parallel Computing

---

## 1. Introduction

Clustering problems have received a widespread concern as one of the most important data mining tasks in various areas. The underlying mechanism of clustering analysis is to partition a dataset into several sub-clusters according to some predefined distance measures [1]. The data samples which belong to the same sub-cluster share the most similar feature values to each other, and the data samples coming to different sub-clusters are the most dissimilar [2]. Numerous attempts have been made to establish a number of clustering algorithms. Considering from the view point of the statistics, the current clustering methods can be simply divided into two categories [2]. The first category methods are not distribution-free approaches such as the expectation and maximization (EM) algorithm

---

\*Corresponding author

Email address: [muyashuang324@126.com](mailto:muyashuang324@126.com) (Yashuang Mu)

11 [3, 4], and the fuzzy  $c$ -directions (FCD) algorithm [5]. The other category methods are  
12 the distribution-free approaches such as K-means algorithm [6] and fuzzy  $c$ -means (FCM)  
13 algorithm [7]. In view of the technical details of realizing clustering process [8], the exist-  
14 ing clustering algorithms can be generally classified into five distinct categories including  
15 partitioning-based clustering algorithms [6, 7, 9, 10, 11, 12], hierarchical-based clustering  
16 algorithms [13, 14, 15, 16, 17], density-based clustering algorithms [18, 19, 20], grid-based  
17 clustering algorithms [21, 22] and model-based clustering algorithms [3, 4, 23, 24]. In ad-  
18 dition, there is also another category of fashionable approaches named ensemble clustering  
19 [25, 26, 27].

20 In partitioning-based clustering algorithms, FCM is a classic and famous clustering  
21 method based on dynamic partition process. Unlike other partitioning-based methods like  
22 K-means [6], FCM is a soft clustering solution which applies fuzzy set theory to describe the  
23 belonging grades of each sample allows the samples belonging to each class with different  
24 grades. As FCM can be easily implemented, and can produce a reasonable result, it is  
25 always a hot topic in data mining fields [28]. However, with quickly increasing data size  
26 involved in many businesses, plenty of traditional clustering models including the famous  
27 FCM algorithm cannot directly deal with large-scale datasets because of some issues like  
28 memory restriction, time complexity and data's complexity [29].

29 The FCM is one of the most powerful and well-known clustering algorithm. It is in-  
30 evitably applying FCM to large-scale data clustering problems. To address the issues that  
31 occur in many traditional machine learning methods, parallel or distributed computing  
32 becomes a frequently-used and trustworthy mechanism [30]. The Map-Reduce based dis-  
33 tributed platforms [31] such as Hadoop [32], Spark [33], etc., provide an effective parallel  
34 computing solution for many data mining operations [34, 35]. Some proposed parallel clus-  
35 tering approaches about FCM can be found in the following literatures. In 2007, [36] pre-  
36 sented a parallel implementation for FCM cluster analysis by using a validity index named  
37 PBM [37] to evaluate the quality of the partition. [38] investigated the parallelization of  
38 FCM algorithm in the framework of Map-Reduce (MR-FCM), where two Map-Reduce jobs  
39 are utilized in each iteration. Furthermore, [39] parallelized a fast fuzzy  $c$ -means method  
40 (MR-FFCM) and focused on underwater image segmentation. The MR-FFCM developed  
41 a two-layer distribution model to group large-scale images and employed an iterative Map-  
42 Reduce process to optimize the clustering process. In 2017, [40] proposed a FCM model  
43 to classify the sentiments of millions of English documents in Cloudera (a parallel network  
44 environment).

45 The above mentioned parallel FCM approaches not only exhibit a good effectiveness  
46 on large-scale data clustering, but also make an extensive and far-reaching influence in  
47 real world applications. Nevertheless, there is a similar characteristic for these methods,  
48 that is, applying the parallel computing technique to the local parts of the standard FCM  
49 algorithm or each iterative process. Obviously, these parallel computing solutions need  
50 to be performed on a real distribution computing platform which is usually constructed  
51 based on a computer cluster. Although the computer cluster has become an efficient

52 method for most engineering and scientific computations, there are some communication  
53 operations among those computers in the cluster. Thus, the communication cost cannot  
54 be avoided in a real distribution computing platform. As the standard FCM algorithm is  
55 an iterative procedure, the previous technique applying parallel computing to local parts  
56 or each iterative process may lead to a lot of communication costs. Motivated by this, in  
57 this study, a simple layered parallel clustering algorithm via FCM technique (LP-FCM)  
58 is developed in the framework of Map-Reduce. Unlike applying parallel computing to  
59 local parts or each iterative process, the proposed LP-FCM algorithm applies Map-Reduce  
60 strategy to the whole of the standard FCM algorithm.

61 In the proposed LP-FCM algorithm, there are only three parallel computing operations  
62 via Map-Reduce model and each operation can be simply treated as a data processing  
63 layer. First, we develop a parallel data partitioning method at the first layer to randomly  
64 partition an original dataset into several subdatasets. Then, the second layer designs a  
65 cluster centers searching method (the processing procedure is somewhat similar to ensemble  
66 clustering approach), where the classic FCM algorithm is employed to search cluster centers  
67 in each Map phase. Each subdataset can produce a set of centers, and these centers are  
68 gathered to a new dataset in Reduce phase. Meanwhile, the classic FCM algorithm is  
69 applied to the new dataset again to confirm the final cluster centers. Finally, based on the  
70 final cluster centers, a parallel data clustering method is organized as the third layer. In  
71 this stage, the original samples are clustered in a parallel manner on the basis of the Euclid  
72 distance between each of the samples and each of the final cluster centers.

73 The major contributions of the designed LP-FCM method can be summarized as two  
74 aspects: (1) the proposed LP-FCM can achieve a good competitive clustering result by  
75 comparing with other famous sequential traditional clustering algorithms; (2) the estab-  
76 lished LP-FCM is a simple and effective method in addressing the memory restriction and  
77 time-consuming problems that arise in large-scale data clustering problems.

78 This study is structured as follows. Section 2 provides some introductions about FCM  
79 algorithm and Map-Reduce. In Section 3, the proposed LP-FCM algorithm is established.  
80 Section 4 shows the experimental results. Finally, the conclusions are made in Section 5.

## 81 **2. FCM algorithm and the framework of Map-Reduce**

82 In this section, we provide a brief review about the basic knowledge of FCM algorithm  
83 and Map-Reduce architecture.

### 84 *2.1. FCM algorithm*

85 FCM algorithm is one of the most popular clustering methods, which realizes the final  
86 partitions through optimizing the basic  $c$ -means objective function. By using the simple  
87 Picard iteration in the first-order conditions for stationary points, the nonlinear minimiza-  
88 tion problem is resolved in FCM algorithm. Bezdek [7] has proven the convergence of FCM  
89 algorithm. The FCM algorithm is subject to the principle that each data point belongs to

90 more than one cluster with different membership values ranging from  $[0, 1]$ . Additionally,  
 91 the sum of the membership values for each data point must be equal to one [2, 38].

92 Let  $X = \{x_i\}_{i=1}^N$  be a collection of instances in  $n$ -dimensional vector space, and  $c$   
 93 ( $2 \leq c \leq N$ ) denotes the number of clusters. An optimal  $c$  partition is realized iteratively  
 94 by minimizing the weighted sum of squared error objective function:

$$J = \sum_{i=1}^N \sum_{j=1}^c (u_{ij})^w d^2(x_i, c_j), \quad (1)$$

95 where  $u_{ij} = u(x_i, c_j) \in [0, 1]$  is the degree of membership of  $x_i$  belonging to the  $j^{\text{th}}$  cluster,  
 96  $w$  ( $1 < w < \infty$ ) is a fuzzification coefficient on each fuzzy membership,  $c_j$  is the center of  
 97 the  $j^{\text{th}}$  cluster, and  $d^2(x_i, c_j)$  is a square distance measure between the instance  $x_i$  and the  
 98 center  $c_j$ . The detailed iterative process can be obtained in Algorithm 1.

---

**Algorithm 1:** The FCM algorithm.

---

**Input:** Weighted exponent of fuzzy membership  $w$ ; threshold  $\varepsilon$  used as a stopping criterion;

dataset  $X = \{x_i\}_{i=1}^N$ .

**Output:** Updated centers  $C$  and fuzzy partition matrix  $U$ .

1 Randomly initialize the fuzzy partition matrix  $U = [u_{ij}(k)]$ , where  $k$  is the iterative times,

$i = 1, 2, \dots, N, j = 1, 2, \dots, c$ .

2 **while** *true* **do**

3     Calculate the cluster centers with  $U^k$ :

$$c_j = \frac{\sum_{i=1}^N u_{ij}^w(k) x_i}{\sum_{i=1}^N u_{ij}^w(k)}. \quad (2)$$

Update the membership matrix  $U^{k+1}$  using:

$$u_{ij}(k+1) = \frac{1}{\sum_{k=1}^c \left( \frac{d(x_i, c_j)}{d(x_i, c_k)} \right)^{\frac{2}{w-1}}}, \quad (3)$$

where  $d(x_i, c_j) = \|x_i - c_j\|^2$ .

4     **if**  $\max_{ij} \|u_{ij}(k+1) - u_{ij}(k)\| < \varepsilon$  **then**

5         **break**.

6     **end**

7 **end**

8 **return**  $C = \{c_1, c_2, \dots, c_c\}$  and  $U = [u_{ij}]$ .

---

## 99 2.2. The framework of Map-Reduce

100 Map-Reduce is a very effective parallel computing model [41, 31]. The simply intent is  
 101 to apply some parallel or distributed algorithm for some large-scale data mining operations  
 102 in a cluster. The cluster may contains several computers or thousands of computers, where  
 103 the Map-Reduce can be realized on different distributed computing platforms to resolve  
 104 various real world problems [38, 42, 43]. There are many Map-Reduce based computing  
 105 platforms. Hadoop is a famous one among this kinds of platforms, which is an open-source  
 106 software for parallel computing or distributed computing with the ability of high reliability,  
 107 scalability and fault tolerance [32].

108 Hadoop can directly provide a very simple programming framework through applying  
 109 Map-Reduce model to deal with large-scale datasets in a distributed computing way. As

110 an open-source software running in a cluster of computers, the primary principle is to  
 111 complete a whole task in the manner of several parallel single subtasks. There are mainly  
 112 two processing phases in a Map-Reduce job: Map and Reduce [44, 45]. The Map phase  
 113 involves a *Mapper* function, and the Reduce phase involves a *Reducer* function, where  
 114 these two functions need to be flexibly implemented by users according to their actual  
 115 situations. Meanwhile, these two functions take some  $\langle key, value \rangle$  pairs as the inputs  
 116 and outputs. A brief illustration about Map-Reduce workflow can be found in Fig. 1.

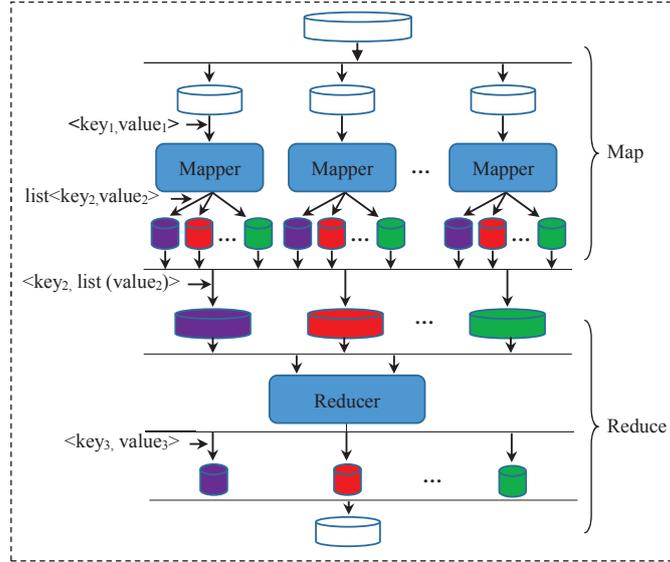


Figure 1: The detailed workflow of Map-Reduce.

117 In general, the *Mapper* function processes the input data and produces some interme-  
 118 diate results in the Mapper phase. The *Mapper* function generates a list of intermediate  
 119  $\langle key, value \rangle$  pairs from a stream of single  $\langle key, value \rangle$  pair. The concrete form can  
 120 be expressed as follows:

$$\text{map } \langle key_1, value_1 \rangle \Longrightarrow \text{list } \langle key_2, value_2 \rangle . \quad (4)$$

121 All the output intermediate  $\langle key, value \rangle$  pairs are grouped by *key* before delivering  
 122 to the *Reducer* function. Then, these intermediate results are fused as the inputs of the  
 123 *Reducer* function during the Reducer phase. The *Reducer* function exports the final results  
 124 (the form is also  $\langle key, value \rangle$  pair) by virtue of aggregation computation. The processing  
 125 procedure can be represented by the following formula:

$$\text{reduce } \langle key_2, \text{list}(value_2) \rangle \Longrightarrow \langle key_3, value_3 \rangle . \quad (5)$$

### 126 3. The layered parallel clustering algorithm

127 We design a layered parallel algorithm via FCM technique, that is LP-FCM, for data  
 128 clustering in this section. The proposed LP-FCM algorithm mainly contains three layers:  
 129 the first layer applies a parallel data partitioning method to randomly divide the original

130 dataset into several subdatasets; the second layer presents a parallel cluster centers search-  
 131 ing method based on the classic FCM algorithm; the third layer uses a parallel clustering  
 132 data method in terms of the final cluster centers. All these three layers are parallel com-  
 133 puting solutions in the framework of Map-Reduce. The main structure of LP-FCM method  
 134 can be illustrated in Fig. 2. In what follows, we will elaborate the layers one by one.

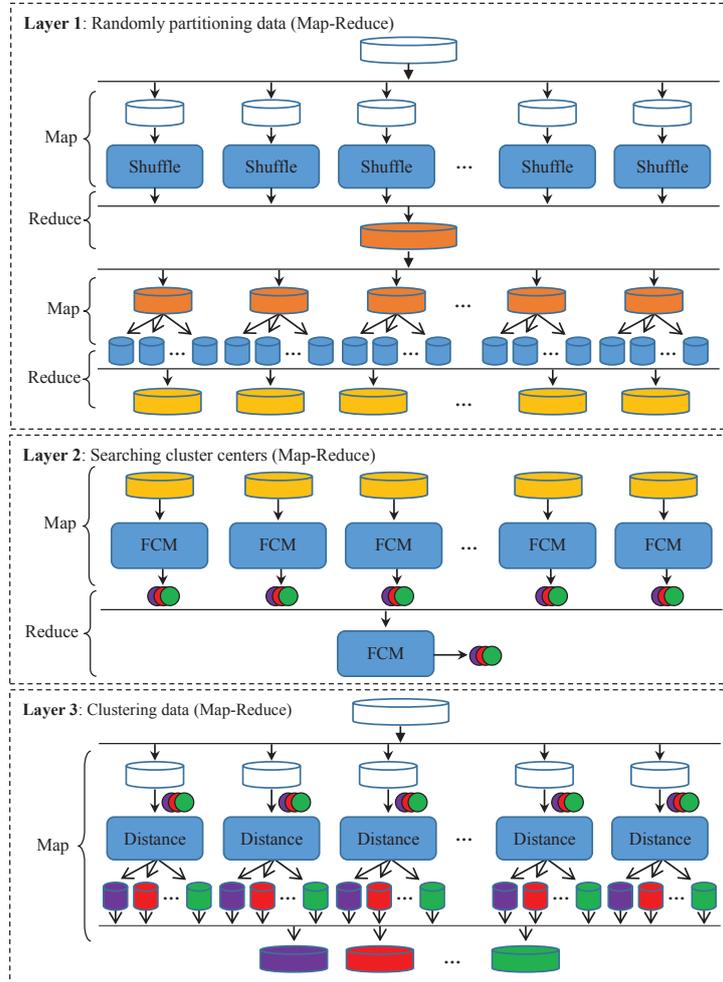


Figure 2: The main structure of LP-FCM algorithm. (1) Layer 1 contains two Map-Reduce jobs. The first job divides the original data (large hollow cylinder) into several sub-datasets (medium hollow cylinders) and shuffles each sub-dataset. Finally, the randomized data (large light red cylinder) are formed. The second job divides the randomized data into several sub-datasets (medium light red cylinders) and split each sub-dataset into several small data sets (small blue cylinders). Finally, the small data sets with same key are combined into one sub-dataset (medium yellow cylinder); (2) Layer 2 applies FCM to each sub-dataset (medium yellow cylinder) to confirm the centers (green purple/red/green small circle). The centers are merged into a new dataset for applying FCM to determine the final centers. (3) Layer 3 labels each sample with the nearest center through calculating the distance between the sample and each center. Finally, the samples (small purple/red/green cylinder) belonging to the same category are gathered to a cluster (medium purple/red/green cylinder).

### 135 3.1. Layer 1: randomly partitioning data

136 In this layer, we develop a parallel data partitioning method to randomly divide a  
 137 dataset into several subdatasets through applying Map-Reduce technique. As shown in

138 Fig. 2, the developed method mainly contains two Map-Reduce jobs. The first Map-  
 139 Reduce job is a parallel data randomizing method, which is used to randomize the data  
 140 samples, and the second Map-Reduce job is a parallel data partitioning method aiming to  
 141 divide the randomized dataset into several subdatasets.

142 Firstly, in order to randomize a large-scale dataset in a parallel processing way, the  
 143 first Map-Reduce job used in this layer considers two parameters in advance: the number  
 144 of used *Mappers* and the number of random values in each *Mapper*. The first parameter  
 145 determines the number of subsets by applying a distributed computing platform to the  
 146 original dataset, and the other parameter determines the number of keys generated in each  
 147 *Mapper*. Suppose the dataset  $X$  can be divided into  $m$  subsets  $X_1, X_2, \dots, X_m$ , the number  
 148 of random values used in each *Mapper* is  $p$ . In the *Mapper* phases, each sample in subset  
 149  $X_j$  will be assigned a random integer value as the output *key*, and the output *value* is the  
 150 sample. Then, the samples with the same *key* are gathered in the *Reducer* phase. Finally,  
 151 depending on the assigned *key*, the samples are stored into some different places. The  
 152 detailed data randomizing process can be found in the following Algorithm 2.

---

**Algorithm 2:** A parallel randomizing data method.

---

**Input:** Dataset  $X$ ; the number of random values  $p$ .  
**Output:** A randomized data  $\tilde{X}$ .

- 1 **Initialize a Hadoop Job *RandomizedDataJob*:**
- 2     Set *RandomizedDataMapper* as the *Mapper* class.
- 3     Set *RandomizedDataReducer* as the *Reducer* class.
- 4     Suppose dataset  $X$  can be partitioned into  $m$  subsets:  $X_1, X_2, \dots, X_m$ .
- 5 **In the  $j^{\text{th}}$  *RandomizedDataMapper*:**
- Input:** Subset  $X_j = \{x_i\}_{i=1}^N$ , where  $x_i$  is the  $i^{\text{th}}$  instance.
- Output:**  $\langle key, value \rangle = \langle n_s, x_i \rangle$ , where  $n_s$  ( $n_s \in [1, p]$ ) is a random integer value.
- 6 **for each  $x_i \in X_j$  do**
- 7     Randomly generate an integer value  $n_s$  ( $n_s \in [1, p]$ ).
- 8     **Mapper Output:**  $\langle key, value \rangle = \langle n_s, x_i \rangle$ .
- 9 **end**
- 10 **In the *RandomizedDataReducer*:**
- Input:**  $\langle key, value \rangle = \langle n_s, list(x_i(j)) \rangle$ , where  $j = 1, 2, \dots, m$ ,  $list(x_i(j))$  denotes the instances with  
    the same key  $n_s$ .
- Output:**  $\langle key, value \rangle = \langle null, list(x_i(j)) \rangle$ .
- 11 Store the instances  $list(x_i(j))$  into  $\tilde{X}$ .
- 12 **Reducer Output:**  $\langle key, value \rangle = \langle null, list(x_i(j)) \rangle$ .
- 13 **return** A randomized data  $\tilde{X}$ .

---

153 Then, the second Map-Reduce job in this layer is aimed to divide the randomized  
 154 dataset into several subdatasets. There are also two parameters that need to be determined  
 155 in advance. The first parameter is the number of used *Mapper*. The second parameter  
 156 is the number of the partitioned subdatasets. If the number of used *Mapper* is  $m$  and  
 157 the number of the partitioned subdatasets is  $\aleph$ , then this Map-Reduce job can divide a  
 158 dataset into  $\aleph$  subdatasets through  $m$  *Mappers*. The following Algorithm 3 illustrates the  
 159 partitioning process in detail.

---

**Algorithm 3:** A parallel partitioning data method.

---

**Input:** Dataset  $X$ ; the number of partitioned subdatasets  $\aleph$ .

**Output:** Subsets  $X^1, X^2, \dots, X^\aleph$ .

1 **Initialize a Hadoop Job *PartitionedDataJob*:**

2     Set *PartitionedDataMapper* as the Mapper class.

3     Set *PartitionedReducer* as the Reducer class.

4     Suppose dataset  $X$  can be partitioned into  $m$  subsets:  $X_1, X_2, \dots, X_m$ .

5 **In the  $j^{\text{th}}$  *PartitionedDataMapper*:**

**Input:** Subset  $X_j = \{x_i\}_{i=1}^N$ , where  $x_i$  is the  $i^{\text{th}}$  instance.

**Output:**  $\langle key, value \rangle = \langle n_s, x \rangle$ , where  $n_s = 1, 2, \dots, \aleph$ .

6 **for each  $x_i \in X_j$  do**

7     Compute the remainder  $n_s$  of  $\frac{i}{|X_j|}$ .

8     **Mapper Output:**  $\langle key, value \rangle = \langle n_s + 1, x_i \rangle$ .

9 **end**

10 **In the *PartitionedDataReducer*:**

**Input:**  $\langle key, value \rangle = \langle n_s, list(x_i(j)) \rangle$ , where  $j = 1, 2, \dots, m$ ,  $list(x_i(j))$  denotes the instances with the same key  $n_s$ .

**Output:**  $\langle key, value \rangle = \langle null, list(x_i(j)) \rangle$ .

11 Store the instances  $list(x_i(j))$  into  $X^{n_s}$ .

12 **Reducer Output:**  $\langle key, value \rangle = \langle null, list(x_i(j)) \rangle$ .

13 **return** Subdatasets  $X^1, X^2, \dots, X^\aleph$ .

---

160 *3.2. Layer 2: searching cluster centers*

161     Suppose dataset  $X$  can be randomly partitioned into  $m$  subdatasets through the first  
162 layer (Layer 1), in Layer 2, a parallel method based on FCM algorithm is established to  
163 search the cluster centers of dataset  $X$  in the framework of Map-Reduce. As there are  $m$   
164 subdatasets, thus the needed number of *Mappers* is  $m$ . If the target number of clusters is  
165  $c$ , then we firstly apply the classic FCM clustering method to each subdataset to confirm  
166 its  $c$  cluster centers in each Mapper phase. All the cluster centers of subdatasets (the total  
167 number of centers is  $m * c$ ) are delivered to the Reducer phase. These centers are gathered  
168 to produce a new small dataset  $\hat{X}$ . Furthermore, the classic FCM algorithm is applied to  
169 the new dataset  $\hat{X}$  again and search the  $c$  cluster centers, where these  $c$  cluster centers are  
170 approximately regarded as the final cluster centers of dataset  $X$ .

171     The following Algorithm 4 shows the searching process of the final cluster centers in  
172 detail. It is easy to find that the data cluster centers of dataset  $X$  are not directly gen-  
173 erated by the classic FCM algorithm, but are produced through applying FCM to several  
174 subdatasets. The classic FCM algorithm, which is treated as a technique, is fused into the  
175 different processing phases of Map-Reduce model. This solution provides a simple way to  
176 endow the classic FCM with the ability to deal with large-scale data clustering problems.

177 *3.3. Layer 3: clustering data*

178     When the cluster centers  $c_1, c_2, \dots, c_c$  are confirmed, the dataset still needs to be parti-  
179 tioned into  $c$  clusters in a parallel computing way following the mechanism of Map-Reduce.  
180 The number of *Mappers* is equal to the number of subsets of dataset  $X$ . In each Mapper

---

**Algorithm 4:** A parallel searching centers method.

---

**Input:** Dataset  $X$ ; the number of clusters  $c$ .

**Output:** The final cluster centers  $c_1, c_2, \dots, c_c$ .

- 1 **Initialize a Hadoop Job *SearchCentersJob*:**
  - 2     Set *SearchCentersMapper* as the *Mapper* class.
  - 3     Set *SearchCentersReducer* as the *Reducer* class.
  - 4     Suppose dataset  $X$  can be partitioned into  $m$  subsets:  $X_1, X_2, \dots, X_m$ .
  - 5 **In the  $j^{\text{th}}$  *SearchCentersMapper*:**
    - Input:** Subset  $X_j = \{x_i\}_{i=1}^N$ , where  $x_i$  is the  $i^{\text{th}}$  instance.
    - Output:**  $\langle key, value \rangle = \langle null, [c_1(X_j), c_2(X_j), \dots, c_c(X_j)] \rangle$ .
  - 6 Apply the classic FCM method (Algorithm 1) on dataset  $X_j$  to confirm cluster centers  $c_1(X_j), c_2(X_j), \dots, c_c(X_j)$ .
  - 7 **Mapper Output:**  $\langle key, value \rangle = \langle null, [c_1(X_j), c_2(X_j), \dots, c_c(X_j)] \rangle$ .
  - 8 **In the *SearchCentersReducer*:**
    - Input:**  $\langle key, value \rangle = \langle null, list([c_1(X_j), c_2(X_j), \dots, c_c(X_j)]) \rangle$ , where  $j = 1, 2, \dots, m$ .
    - Output:**  $\langle key, value \rangle = \langle null, [c_1, c_2, \dots, c_c] \rangle$ .
  - 9 Gather all centers and form a new dataset  $\hat{X}$ .
  - 10 Apply the classic FCM method (Algorithm 1) on dataset  $\hat{X}$  to confirm cluster centers  $c_1, c_2, \dots, c_c$ .
  - 11 **Reducer Output:**  $\langle key, value \rangle = \langle null, [c_1, c_2, \dots, c_c] \rangle$ .
  - 12 **return**  $c_1, c_2, \dots, c_c$ .
- 

181 phase, the Euclidean distance between each sample in  $X_j$  and each center is firstly calcu-  
182 lated. Then, the index of centers with the nearest distance to a sample is selected as the  
183 *key*, meanwhile, the sample is treated as the *value*. As the outputs with the same *key* of  
184 each *Mapper* task can be stored into the same place, which is similar to the Reduce phase,  
185 for this reason, only the Mapper phase is utilized in this layer. The details can be found  
186 in the following Algorithm 5.

---

**Algorithm 5:** A parallel clustering data method.

---

**Input:** Dataset  $X$ ; the final cluster centers  $c_1, c_2, \dots, c_c$ .

**Output:** The final clusters  $X_1, X_2, \dots, X_c$ .

- 1 **Initialize a Hadoop Job *ClusterDataJob*:**
  - 2     Set *ClusterDataMapper* as the *Mapper* class.
  - 3     Set *ClusterDataReducer* as the *Reducer* class.
  - 4     Suppose dataset  $X$  can be partitioned into  $m$  subsets:  $X_1, X_2, \dots, X_m$ .
  - 5 **In the  $j^{\text{th}}$  *ClusterDataMapper*:**
    - Input:** Subset  $X_j = \{x_i\}_{i=1}^N$ , where  $x_i$  is the  $i^{\text{th}}$  instance.
    - Output:**  $\langle key, value \rangle = \langle null, [X_1, X_2, \dots, X_c] \rangle$ .
  - 6 **for each**  $x_i \in X_j$  **do**
    - 7     Compute Euclidean distance  $d(x_i, c_k)$ , where  $k = 1, 2, \dots, c$ .
    - 8     Select  $k^* = \arg \min_k d(x_i, c_k)$  and store  $x_i$  to  $X_{k^*}$ .
    - 9     **Mapper Output:**  $\langle key, value \rangle = \langle k^*, x_i \rangle$ .
  - 10 **end**
  - 11 **return**  $X_1, X_2, \dots, X_c$ .
-

187 *3.4. The proposed LP-FCM clustering algorithm*

188 On basis of the previous work, the layered parallel clustering algorithm via FCM tech-  
189 nique (LP-FCM) is established in this section. The detailed pseudocode is given in Algo-  
190 rithm 6.

---

**Algorithm 6:** The proposed LP-FCM algorithm.

---

**Input:** Dataset  $X$ ; the number of centers  $c$ .

**Output:** The final clusters  $X_1, X_2, \dots, X_c$ .

- 1 Randomize the dataset  $X$  by Algorithm 2 (Layer 1).
  - 2 Partition  $X$  into  $m$  subdatasets  $\{X_j\}_{j=1}^m$  by Algorithm 3 (Layer 1).
  - 3 Search the final cluster centers  $c_1, c_2, \dots, c_c$  by Algorithm 4 (Layer 2).
  - 4 Divide dataset  $X$  into  $c$  clusters  $X_1, X_2, \dots, X_c$  by Algorithm 5 (Layer 3).
  - 5 **return**  $X_1, X_2, \dots, X_c$ .
- 

191 As shown in Algorithm 6, there are three layers or four steps in the pseudocode. First,  
192 the dataset  $X$  is randomized and partitioned by Algorithm 2 and Algorithm 3, respectively.  
193 Then, we employ Algorithm 4 to search the final cluster centers. Finally, the dataset  $X$   
194 is divided into  $c$  clusters  $X_1, X_2, \dots, X_c$  by Algorithm 5. All these steps are some parallel  
195 implementations in the framework of Map-Reduce. The proposed LP-FCM approach incor-  
196 porates the classical FCM method to the procedure of Map-Reduce in a simple and specific  
197 way, which can effectively resolve some issues in dealing with large-scale data clustering  
198 for the traditional FCM algorithm.

199 *3.5. Analysis of the time and space complexity*

200 Suppose there are  $n$  samples and each sample has  $f$  attributes, the number of iterations  
201 is  $k$ , the number of clusters is  $c$ , it can be analyzed that the time and space complexity  
202 of the traditional FCM algorithm are  $O(k * c * n * f)$  and  $O(n * f)$ , respectively. For  
203 the proposed LP-FCM algorithm, the time and space complexity is mainly related to the  
204 two phases of Map-Reduce. If the original dataset can be partitioned into  $m$  subsets in  
205 Map phase, then the time complexity is  $O(\frac{k * c * n * f}{m})$  and the space complexity is  $O(\frac{n * f}{m})$ .  
206 In the Reduce phase, as the number of gathered centers is  $c * m$ , the time complexity  
207 is  $O(k * c * c * m * f)$  and the space complexity is  $O(k * m * f)$ . Thus, the total time  
208 and space complexity of the proposed LP-FCM are  $O(\frac{k * c * n * f}{m}) + O(k * c * c * m * f)$   
209 and  $O(\frac{n * f}{m}) + O(k * m * f)$ , respectively. However, in practice, due to the reasons like a  
210 cluster system with heterogeneous hardware, slightly different data partition sizes, network  
211 latencies, etc, there are usually some communication costs in a cluster of computers, which  
212 is a challenging problem in parallel computing or distributed computing system.

213 **4. Experimental studies**

214 To assess the performance of the proposed LP-FCM approach, several numerical ex-  
215 perimental studies are described. First, Section 4.1 covers some descriptions about the  
216 computing environment and the datasets used in the experiments. Then, in order to verify

217 the feasibility of the proposed LP-FCM algorithm, a comparative analysis on clustering  
 218 accuracy with some famous random initialization clustering methods is offered in Section  
 219 4.2. Furthermore, Section 4.3 and Section 4.5 include some analysis about the parallelism  
 220 in terms of the execution time and some evaluation indices on some generated datasets.

#### 221 4.1. Computing environment and data sets

222 In this study, we implement the proposed LP-FCM algorithm with Java language based  
 223 on the open source Hadoop distribution computing platform [32] (the version of software is  
 224 2.6.0), meanwhile, some data structures in WEKA [46] (the version of software is 3.6.9) are  
 225 also utilized. All sequential methods are implemented on base of the *statistics and machine*  
 226 *learning* toolbox in the platform of MATLAB (the version of software is R2015b). All the  
 227 experiments related to the sequential modes are carried out on a machine with Intel Core  
 228 i5-4590 3.30 GHZ, 8 GB RAM and Ubuntu 14.04.1 LTS (64 Bit) OS, and all experiments  
 229 related to parallel modes are conducted on a small distribution computing cluster, where it  
 230 mainly contains one host machine with Intel Core i5-4590 3.30 GHZ, 8 GB RAM, Ubuntu  
 231 14.04.1 LTS (64 Bit) OS and six servant machines with Intel Core i3 2.93 GHZ, 4 GB  
 232 RAM, Ubuntu 14.04.1 LTS (64 Bit) OS.

233 The experiments totally employ 21 UCI benchmark datasets [47] to evaluate the se-  
 234 quential clustering algorithms and the proposed parallel clustering algorithm in this study.  
 235 Specially, the samples with missing values are deleted from the original data sets. The de-  
 236 tailed compositions related to these datasets are exhibited in Table 1. The first 20 datasets  
 237 are mainly used to compare with some sequential clustering algorithms for the feasibility  
 238 analysis, and the last dataset is aimed to test the parallelism.

Table 1: Compositions of the benchmark data sets.

No.	Data sets	Attributes	Class	Instances
1	Banknote	4	2	1,372
2	Blood	4	2	748
3	Breast-cancer-w-o	9	2	699
4	Breast-cancer-w-d	30	2	569
5	Ecoli	7	8	366
6	Fertility	9	2	100
7	Glass	9	7	214
8	Hayes-roth	4	3	160
9	Heart-statlog	13	2	270
10	Iris	4	3	150
11	Lenses	4	3	24
12	Magic	10	2	19,020
13	Messidor	19	2	1,151
14	Skin	3	2	245,057
15	Seeds	7	3	210
16	Segment	19	2	2,310
17	Sonar	60	2	208
18	Waveform	40	3	5,000
19	Wine	13	3	178
20	Wilt	5	2	4,839
21	Coverttype	54	7	581,012

239 4.2. Comparisons with K-means, K-medoids, FCM and MinMaxK-means

240 In order to test the feasibility of the proposed LP-FCM algorithm, we compare it with  
 241 several famous random initialization clustering methods: K-means [6], K-medoids [9], FCM  
 242 [7] and MinMaxK-means [48]. The classic K-means, K-medoids and FCM are implemented  
 243 by using the related functions in MATLAB toolboxes, where the parameter ‘Start’ is set  
 244 as ‘sample’ (if exists) and others parameters are set as their defaults. The implementation  
 245 of MinMaxK-means is also developed by Matlab codes, in which the parameters are used  
 246 as their defaults.

247 During the experiments, all these clustering methods are executed 50 times for each  
 248 benchmark dataset, and the real number of class labels is taken as the initial number of  
 249 clusters. The clustering accuracies [28, 49] on all benchmark datasets are recorded in Table  
 250 2, where the values are the means values of 50 independent runs. Meanwhile, the number  
 251 of processors used in the proposed LP-FCM approach is set as 3. The results marked  
 252 with bold face denote that the current method is the best one among all methods on the  
 253 corresponding dataset.

Table 2: The comparison of clustering accuracies on 20 data sets.

Data sets	K-means	K-medoids	FCM	MinMaxK-means	LP-FCM
Banknote	<b>0.6122</b>	0.5951	0.6093	0.6028	0.5994
Blood	0.7247	0.7259	0.7072	<b>0.7262</b>	0.7005
Breast-cancer-w-o	<b>0.9610</b>	0.9590	0.9561	0.9213	0.9458
Breast-cancer-w-d	0.8541	0.8541	0.8541	0.8313	<b>0.8660</b>
Ecoli	0.5397	0.5932	0.5107	0.5483	<b>0.6393</b>
Fertility	0.6064	0.5312	0.5514	0.5876	<b>0.6160</b>
Glass	0.4995	<b>0.5349</b>	0.4907	0.5015	0.5083
Hayes-roth	<b>0.4648</b>	0.4294	0.3453	0.4606	0.4613
Heart-statlog	0.5899	<b>0.6074</b>	0.5926	0.5926	0.5984
Iris	0.8340	0.8939	0.8933	0.8867	<b>0.9335</b>
Lenses	0.4825	0.4975	0.4167	0.4725	<b>0.4992</b>
Magic	<b>0.6491</b>	0.6168	0.5780	0.6392	0.5969
Messidor	0.5065	0.5155	0.5352	0.5039	<b>0.5585</b>
Skin	0.5513	<b>0.5514</b>	0.5511	0.5399	0.5382
Seeds	0.8930	<b>0.8952</b>	<b>0.8952</b>	<b>0.8952</b>	0.8734
Segment	<b>0.6956</b>	0.6935	0.6420	0.6645	0.5695
Sonar	0.5450	0.5203	0.5529	0.5481	<b>0.5704</b>
Waveform	0.5142	<b>0.5768</b>	0.5248	0.5162	0.5747
Wine	0.6637	0.6820	0.6854	0.6685	<b>0.6942</b>
Wilt	0.5200	0.5211	0.5015	<b>0.5226</b>	0.5012
Avg.	0.6354	0.6397	0.6197	0.6315	<b>0.6426</b>
	(5/20)	(5/20)	(1/22)	(3/20)	<b>(8/20)</b>

254 From the values recorded in Table 2, we can intuitively find that, among these al-  
 255 gorithms, the classic K-means, K-medoids, FCM, and MinMaxKmeans methods achieve  
 256 better performance on 5, 5, 1 and 3 datasets among 20 datasets, while the proposed LP-  
 257 FCM algorithm can achieve 8 datasets among 20 datasets. According to the average of  
 258 clustering accuracies on all these benchmark datasets, the proposed LP-FCM approach  
 259 also obtains a similar result with other clustering algorithms.

260 Furthermore, in order to check whether there is a significant difference among these

261 clustering methods, we provide a statistical discussion on these computational results by  
 262 employing Friedman test [50]. In the Friedman test, the null-hypothesis is that the tested  
 263 indices are equivalent. At the significance level  $\alpha = 0.05$ , the critical value with 4 and  
 264 76 degrees of freedom is 2.4920, but the  $F_F$  derived from the Table 2 is 1.2748, which  
 265 is smaller than 2.4920. Thus, there is not a significant difference among these clustering  
 266 methods.

267 From the above discussions, the results deduce that the proposed LP-FCM method  
 268 can keep the same level of performance on clustering accuracy with some classic clustering  
 269 methods, which validates the feasibility of the proposed LP-FCM algorithm.

#### 270 4.3. The execution time analysis

271 As the proposed LP-FCM algorithm is a parallel solution for clustering problems, it is  
 272 necessary to study the execution time. In this experiment, we examine the execution time  
 273 on *Covertypes* dataset. The dataset contains 54 conditional attributes and one decision  
 274 attribute with 7 categories. The original *Covertypes* dataset includes 581,012 instances,  
 275 however, to be able to clearly present the parallelism, we generate 6 datasets from the  
 276 original dataset by bootstrap technique [51] to test the execution time for data clustering.  
 277 In Table 3, the generated 6 datasets are described in detail.

Table 3: The detailed information of the generated datasets.

No.	Data sets	Attributes	Class	Instances
1	Covertypes (25MB)	54	7	202,621
2	Covertypes (50MB)	54	7	405,242
3	Covertypes (75MB)	54	7	607,863
4	Covertypes (100MB)	54	7	810,484
5	Covertypes (125MB)	54	7	1,013,105
6	Covertypes (150MB)	54	7	1,215,726

278 The execution time of Algorithm 6 at six different processors is summarized in Table 4,  
 279 where the results with mark “—” denote that there is a memory warning in the experiment.  
 280 According to the values, it is obvious that, in most cases, the execution time of LP-  
 281 FCM algorithm presents a downtrend with the increasing number of the processors, which  
 282 reflects the efficiency of parallel solution in reducing execution time. Especially, for the  
 283 case of *Covertypes (25MB)*, the value with six processes is larger than the value with five  
 284 processes, which is unnatural comparing with other cases. The reason is that there is a  
 285 communication time existing among these processors. When the amount of data is smaller  
 286 and the number of processors is bigger, the communication time takes a large proportion  
 287 in the total execution time. In additional, when the amount of data is increased, the less  
 288 processors cannot deal with the dataset because of some memory restrictions, but this  
 289 issue can be resolved through utilizing more processors. Taking the case of *Covertypes*  
 290 *(100MB)*, when two processors are employed to deal with the dataset, a memory warning  
 291 is appeared. However, we can observe that the memory warning is disappeared when the  
 292 number of processors is raised to three or more.

Table 4: The execution time (seconds) of LP-FCM.

Data sets	No. of processors					
	One	Two	Three	Four	Five	Six
Covertypes (25MB)	947.278	530.862	392.053	343.638	323.866	346.952
Covertypes (50MB)	—	1746.253	1020.711	788.154	672.986	568.792
Covertypes (75MB)	—	—	1503.130	1302.679	921.933	805.442
Covertypes (100MB)	—	—	2311.940	1735.194	1431.983	1193.799
Covertypes (125MB)	—	—	—	2454.673	2037.295	1485.984
Covertypes (150MB)	—	—	—	—	2559.292	1776.333

#### 293 4.4. The comparison on reducing time between LP-FCM and MR-FCM

294 In order to further present the effectiveness of the proposed LP-FCM algorithm on  
 295 reducing execution time, we make a comparison on execution time between LP-FCM and  
 296 the parallel clustering algorithm MR-FCM [38]. The main similarity of the two methods  
 297 is that both of them are the parallelization techniques of the classic FCM algorithm. The  
 298 main difference is that the proposed LP-FCM algorithm applies parallel computing to the  
 299 whole of the classic FCM approach, while the parallel MR-FCM algorithm applies parallel  
 300 computing to local parts or each iterative process. The used datasets in the experiments  
 301 are selected from Table 3.

302 Firstly, we fix the number of processors and employ the datasets *Covertypes (50MB)*,  
 303 *Covertypes (50MB)*, *Covertypes (75MB)*, *Covertypes (100MB)* in Table 3 to test the execution  
 304 time. During the experiments, both of the two algorithms employ the same parameters,  
 305 and are executed on the same cluster where the number of processors are set as 3. Through  
 306 conducting the experiments under the same conditions, the values are obtained and dis-  
 307 played in Fig. 3. We can easily find that the proposed LP-FCM algorithm can evidently  
 308 save the execution time comparing with the parallel MR-FCM algorithm.

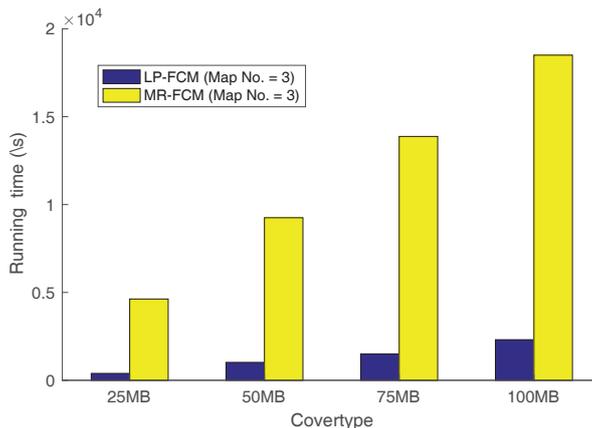


Figure 3: The execution time on different datasets.

309 Then, we fix the data size (where the used dataset is *Covertypes (50MB)*) and adjust the  
 310 number of processors from 2 to 5. From the comparison of results as shown in Fig. 4, two  
 311 conclusions can be easily summarized. (1) The execution time presents a downward trend  
 312 with the increase in the number of processors. (2) The proposed LP-FCM parallel solution

313 has more advantages on execution time than the parallel MR-FCM clustering algorithm.

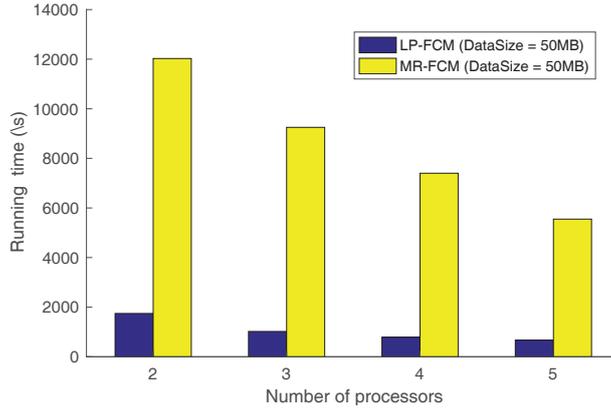


Figure 4: The execution time on different processors.

#### 314 4.5. The parallel performance analysis

315 In addition, we also compute some evaluation indices such as Speedup, Scaleup and  
 316 Sizeup [52] to further verify the parallel performance of the proposed LP-FCM algorithm.  
 317 Speedup represents how much a parallel system with  $m$  processors is faster in time than a  
 318 system with one processor. Scaleup means the ability of a parallel system with  $m$  processors  
 319 to process  $m$  times data during the same original computing time. Sizeup measures how  
 320 much longer time is consumed while dealing with  $m$  times original data than dealing with  
 321 the original data in a given parallel system.

322 Here, we treat the dataset *Coverttype (25MB)* as the original dataset in Table 3, and the  
 323 datasets *Coverttype (50MB)*, *Coverttype (75MB)*, *Coverttype (100MB)*, *Coverttype (125MB)*  
 324 and *Coverttype (150MB)* are respectively regarded as the corresponding times datasets of  
 325 *Coverttype (25MB)*. The execution time has been recorded in Table 4. The definitions of  
 326 Speedup and Sizeup require that the missing values must be known. Thus, in order to  
 327 calculate smoothly, we only use the known values for these evaluation indices. In addition, as  
 328 the reasons such as the cluster system with heterogeneous hardware, slightly different data  
 329 partition sizes, network latencies, etc, some communication costs exist in the cluster of  
 330 computers. In what follows, we analyze the parallel evaluation indices in detail.

331 Firstly, we compute the Speedup and present it in Fig. 5. The number of processors  
 332 is varied from 1 to 6, and the original dataset is *Coverttype (25MB)*. As there are some  
 333 missing values for some datasets, thus we only use dataset *Coverttype (25MB)* to depict the  
 334 performance of Speedup. From the tendency of the curve, we can obviously find that the  
 335 Speedup is highly correlated to the number of processors. With the increase in number of  
 336 processors, the Speedup is gradually increasing. However, the Speedup does not exhibit  
 337 an increasing tendency. When the number of processors is five, the Speedup reaches the  
 338 maximum. That reason is that there are more communication costs when the number of  
 339 processors increases, and the communication costs take a bigger proportion in the total  
 340 execution time.

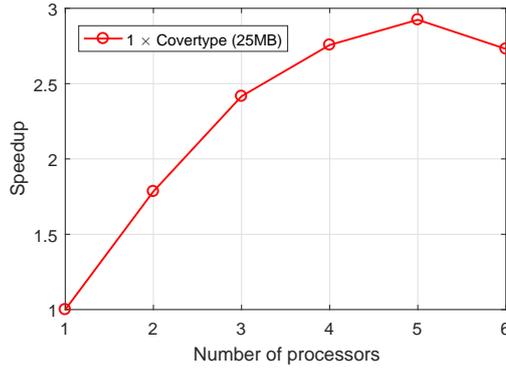


Figure 5: Speedup of the proposed algorithm for real datasets.

341 Then, for the Scaleup performance, the following Fig. 6 depicts the tendency. According  
 342 to the definition, the used datasets are 1 to 6 times of original data *Covertypes* (25MB), and  
 343 the number of processors is equal to the times for each case. In an ideal parallel system,  
 344 Scaleup should be equal to 1. However, as the reason of communication cost among the  
 345 processors, the Scaleup exhibits a downward tendency. Meanwhile, when the data and the  
 346 processors increase, the Scaleup reduces slowly, or shows a fluctuation trend.

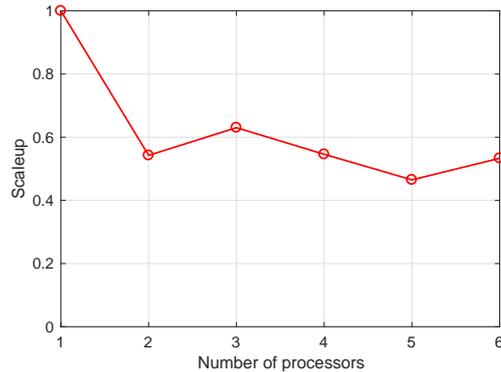


Figure 6: Scaleup of the proposed algorithm for real datasets.

347 Furthermore, we also study the Sizeup performance of the proposed LP-FCM algorithm.  
 348 During the experiment, the processors are fixed as 5 and 6, respectively. Meanwhile, the  
 349 used datasets are 1 to 6 times of the original data *Covertypes* (25MB), respectively. The  
 350 detailed tendencies are summarized in Fig. 7. From the curves, the proposed approach  
 351 shows a good Sizeup performance in the given parallel system.

## 352 5. Conclusions

353 In this study, a parallel algorithm named LP-FCM is proposed for data clustering prob-  
 354 lems. The LP-FCM clustering approach mainly contains three layers: randomly partition-  
 355 ing data, searching cluster centers and clustering data. All the layers are implemented in  
 356 the framework of Map-Reduce. Especially, in the second layer, the famous FCM technique  
 357 is fused to the different phases of Map-Reduce model. During the experimental studies, the

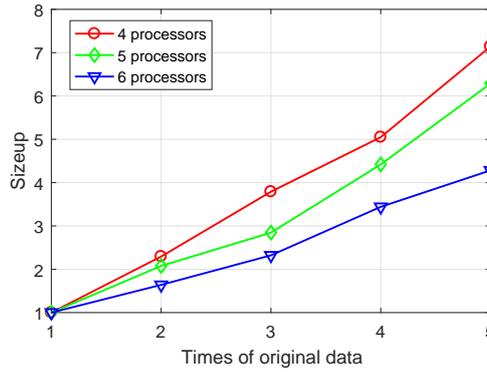


Figure 7: Sizeup of the proposed algorithm for real datasets.

358 feasibility is evaluated by comparing with some traditional clustering algorithms, and the  
 359 parallelism is examined through analyzing the execution time and some evaluation indices  
 360 for parallel performance.

361 At last, it deserves to point out three aspects of this paper: (1) there are three layers in  
 362 the proposed approach, where each layer can be treated as a processing step for other large-  
 363 scale data mining algorithms. For example, the first layer can be regarded as a parallel data  
 364 pre-processing step, and the second layer can be applied to confirm the parameters of fuzzy  
 365 sets in a large-scale fuzzy system; (2) there are also other kinds of clustering algorithms  
 366 such as density-based methods. The idea of effectively parallelize other cluster algorithms  
 367 is also a meaningful theme.

### 368 Acknowledgments

369 This work is supported by the Research Foundation for Advanced Talents (2019BS007,  
 370 31401204) of Henan University of Technology. This work is supported by the Research  
 371 Foundation for Advanced Talents of Henan University of Technology (No. 2019BS007),  
 372 the Open Fund of Key Laboratory of Grain Information Processing and Control (Henan  
 373 University of Technology), Ministry of Education (No. KFJJ-2020-112), and the National  
 374 Natural Science Foundation of China under Grants (No. 62006071).

### 375 Ethical approval

376 The authors declare that there is no conflict of interests regarding the publication of  
 377 this paper.

### 378 Funding

379 This study is not supported by any source or any organizations.

### 380 Conflict of interest

381 The authors declare that there is no conflict of interests regarding the publication of  
 382 this paper.

## 383 Informed Consent

384 Informed Consent was not required as no human or animals were involved.

## 385 Authors' Contributions

386 Yashuang Mu: writing-original draft and model visualization; Wei Wei: methodology  
387 and data analyses; Hongyue Guo: writing-review and data analyses; Lijun Sun: writing-  
388 review and supervision.

## 389 References

- 390 [1] E. Zhu and R. Ma, "An effective partitional clustering algorithm based on new clus-  
391 tering validity index," *Applied Soft Computing*, vol. 71, no. C, pp. 608–621, 2018.
- 392 [2] O. Kesemen, Ö. Tezel, and E. Özkul, "Fuzzy c-means clustering algorithm for direc-  
393 tional data (FCM4dd)," *Expert Systems with Applications*, vol. 58, pp. 76–82, 2016.
- 394 [3] A. P. Dempster, N. M. Laird, and D. B. Rubin, "Maximum likelihood from incomplete  
395 data via the EM algorithm," *Journal of the Royal Statistical Society*, vol. 39, no. 1,  
396 pp. 1–38, 1977.
- 397 [4] G. J. McLachlan and K. E. Basford, "Mixture models: Inference and applications to  
398 clustering," *Inference & Applications to Clustering*, vol. 38, no. 2, 1988.
- 399 [5] M. S. Yang and J. A. Pan, "On fuzzy clustering of directional data," *Fuzzy Sets and*  
400 *Systems*, vol. 91, no. 3, pp. 319–326, 1997.
- 401 [6] J. Macqueen, "Some methods for classification and analysis of multivariate observa-  
402 tions," in *Proc. of Berkeley Symposium on Mathematical Statistics and Probability*,  
403 1965, pp. 281–297.
- 404 [7] J. C. Bezdek, R. Ehrlich, and W. Full, "FCM: The fuzzy c-means clustering algorithm,"  
405 *Computers & Geosciences*, vol. 10, no. 2, pp. 191–203, 1984.
- 406 [8] A. Fahad, N. Alshatri, Z. Tari, and A. Alamri, "A survey of clustering algorithms for  
407 big data: Taxonomy and empirical analysis," *Emerging Topics in Computing IEEE*  
408 *Transactions on*, vol. 2, no. 3, pp. 267–279, 2014.
- 409 [9] H. S. Park and C. H. Jun, "A simple and fast algorithm for K-medoids clustering,"  
410 *Expert Systems with Applications*, vol. 36, no. 2, pp. 3336–3341, 2009.
- 411 [10] R. T. Ng and J. Han, "Efficient and effective clustering methods for spatial data  
412 mining," in *VLDB Conference*, 1994, pp. 144–155.
- 413 [11] L. Kaufman and P. J. Rousseeuw, *Finding Groups in Data: An Introduction to Cluster*  
414 *Analysis*. DBLP, 1990.

- 415 [12] Ng, T. Raymond, Han, and Jiawei, “CLARANS: A method for clustering objects for  
416 spatial data mining,” *IEEE Transactions on Knowledge & Data Engineering*, vol. 14,  
417 no. 5, pp. 1003–1016, 2002.
- 418 [13] T. Zhang, R. Ramakrishnan, and M. Livny, “BIRCH: An efficient data clustering  
419 method for very large databases,” in *Proceedings of the 1996 ACM SIGMOD Inter-  
420 national Conference on Management of Data*, vol. 25, no. 2, 1994, pp. 103–114.
- 421 [14] S. Guha, R. Rastogi, and K. Shim, “CURE: An efficient clustering algorithm for large  
422 databases,” *Information Systems*, vol. 26, no. 1, pp. 35–58, 1998.
- 423 [15] S. Guha, R. Rastogi, and K. Shim, “ROCK: A robust clustering algorithm for cate-  
424 gorical attributes,” *Information Systems*, vol. 25, no. 5, pp. 345–366, 1999.
- 425 [16] G. Karypis, E. H. Han, and V. Kumar, *Chameleon: Hierarchical Clustering Using  
426 Dynamic Modeling*. IEEE Computer Society Press, 1999.
- 427 [17] A. N. Mahmood, C. Leckie, and P. Udaya, “An efficient clustering scheme to exploit  
428 hierarchical data in network traffic analysis,” *IEEE Transactions on Knowledge &  
429 Data Engineering*, vol. 20, no. 6, pp. 752–767, 2008.
- 430 [18] M. Ester, H. P. Kriegel, and X. Xu, “A density-based algorithm for discovering clus-  
431 ters a density-based algorithm for discovering clusters in large spatial databases with  
432 noise,” in *International Conference on Knowledge Discovery and Data Mining*, 1996,  
433 pp. 226–231.
- 434 [19] M. Ankerst, M. M. Breunig, H. P. Kriegel, and J. Sander, “OPTICS: Ordering points  
435 to identify the clustering structure,” in *International Conference on Management of  
436 Data and Symposium on Principles of Database Systems*, 1999, pp. 49–60.
- 437 [20] X. Xu, M. Ester, and H. P. Kriegel, “A distribution-based clustering algorithm for  
438 mining in large spatial databases,” in *International Conference on Data Engineering,  
439 1998. Proceedings*, 1998, pp. 324–331.
- 440 [21] G. Sheikholeslami, S. Chatterjee, and A. Zhang, “Wavecluster: A multi-resolution  
441 clustering approach for very large spatial databases,” in *International Conference on  
442 Very Large Data Bases*, 1998, pp. 428–439.
- 443 [22] W. Wang, J. Yang, and R. R. Muntz, “STING: A statistical information grid approach  
444 to spatial data mining,” pp. 186–195, 1997.
- 445 [23] Fisher and H. Douglas, “Knowledge acquisition via incremental conceptual clustering,”  
446 *Machine Learning*, vol. 2, no. 2, pp. 139–172, 1987.
- 447 [24] J. H. Gennari, P. Langley, and D. Fisher, “Models of incremental concept formation,”  
448 *Artificial Intelligence*, vol. 40, no. 1, pp. 11–61, 1989.
- 449 [25] A. Strehl, J. Ghosh, “Cluster ensembles: A knowledge reuse framework for combining  
450 partitionings,” *Journal of Machine Learning Research*, vol. 3, no. 3, pp. 583–617, 2002.

- 451 [26] S. Vega-Pons, J. Ruiz-Shulcloper, “A survey of clustering ensemble algorithms,” *Journal of Pattern Recognition and Artificial Intelligence*, vol. 25, no. 03, pp. 337–372,  
452 2011.  
453
- 454 [27] M. Mojarad, S. Nejatian, H. Parvin, and M. Majid, “A fuzzy clustering ensemble  
455 based on cluster clustering and iterative Fusion of base clusters,” *Applied Intelligence*,  
456 vol. 47, no. 7, pp. 2567–2581, 2019.
- 457 [28] H. X. Pei, Z. R. Zheng, C. Wang, C. N. Li, and Y. H. Shao, “D-FCM: Density based  
458 fuzzy c-means clustering algorithm with application in medical image segmentation,”  
459 *Procedia Computer Science*, vol. 122, pp. 407–414, 2017.
- 460 [29] Y. Mu, X. Liu, and L. Wang, “A pearson’s correlation coefficient based decision tree  
461 and its parallel implementation,” *Information Sciences*, vol. 435, pp. 40–58, 2018.
- 462 [30] Y. Mu, L. Wang, and X. Liu, “A fast rank mutual information based decision tree  
463 and its implementation via map-reduce,” *Concurrency and Computation: Practice  
464 and Experience*, vol. 30, no. 10, 2018.
- 465 [31] J. Dean and S. Ghemawat, “MapReduce: Simplified data processing on large clusters,”  
466 *Communications of the ACM*, vol. 51, no. 1, pp. 107–113, 2008.
- 467 [32] Apache Hadoop. Available online: <http://hadoop.apache.org/>.
- 468 [33] Apache Spark. Available online: <https://spark.apache.org/>.
- 469 [34] J. Dean and S. Ghemawat, “MapReduce: A flexible data processing tool,” *Communi-  
470 cations of the ACM*, vol. 53, no. 1, pp. 72–77, 2010.
- 471 [35] L. Zhang, K. Li, Y. Xu, J. Mei, F. Zhang, and K. Li, “Maximizing reliability with en-  
472 ergy conservation for parallel task scheduling in a heterogeneous cluster,” *Information  
473 Sciences*, vol. 319, no. C, pp. 113–131, 2015.
- 474 [36] M. V. Modenesi, M. C. A. Costa, A. G. Evsukoff, and N. F. F. Ebecken, “Parallel  
475 fuzzy c-means cluster analysis,” in *International Conference on High Performance  
476 Computing for Computational Science*, pp. 52–65.
- 477 [37] M. K. Pakhira, S. Bandyopadhyay, and U. Maulik, “Validity index for crisp and fuzzy  
478 clusters,” *Pattern Recognition*, vol. 37, no. 3, pp. 487–501, 2004.
- 479 [38] S. A. Ludwig, “MapReduce-based fuzzy c-means clustering algorithm: implementation  
480 and scalability,” *Journal of Machine Learning and Cybernetics*, vol. 6, no. 6, pp. 923–  
481 934, 2015.
- 482 [39] L. Xiu, J. Song, Z. Fan, X. Ouyang, and S. U. Khan, “MapReduce-based fast fuzzy  
483 c-means algorithm for large-scale underwater image segmentation,” *Future Generation  
484 Computer Systems*, vol. 65, no. C, pp. 90–101, 2016.

- 485 [40] V. N. Phu, N. D. Dat, V. T. N. Tran, V. T. N. Chau, and T. A. Nguyen, “Fuzzy  
486 C-means for english sentiment classification in a distributed system,” *Applied Intelli-*  
487 *gence*, vol. 46, no. 3, pp. 717–738, 2017.
- 488 [41] S. Shankland, “Google spotlights data center inner workings,” *CNet News Blog*, *posted*  
489 *on May*, vol. 30, p. 2008, 2008.
- 490 [42] J. Chen, K. Li, Z. Tang, K. Bilal, S. Yu, C. Weng, and K. Li, “A parallel random forest  
491 algorithm for big data in a spark cloud computing environment,” *IEEE Transactions*  
492 *on Parallel & Distributed Systems*, vol. 28, no. 4, pp. 919–33, 2017.
- 493 [43] G. Tang, W. Yang, K. Li, Y. Ye, G. Xiao, and K. Li, “An iteration-based hybrid  
494 parallel algorithm for tridiagonal systems of equations on multi-core architectures,”  
495 *Concurrency and Computation: Practice and Experience*, vol. 27, no. 17, pp. 5076–  
496 5095, 2015.
- 497 [44] T. White, *Hadoop: The Definitive Guide, 4th Edition*. O’Reilly Media, 2015.
- 498 [45] I. Triguero, D. Peralta, J. Bacardit, S. Garca, and F. Herrera, “MRPR: A MapReduce  
499 solution for prototype reduction in big data classification,” *Neurocomputing*, vol. 150,  
500 no. 150, pp. 331–345, 2015.
- 501 [46] M. Hall, E. Frank, G. Holmes, B. Pfahringer, P. Reutemann, and I. H. Witten, “The  
502 WEKA data mining software: An update,” *ACM SIGKDD Explorations Newsletter*,  
503 vol. 11, no. 1, pp. 10–18, 2009.
- 504 [47] D. Dheeru and E. K. Taniskidou, 2017, UCI Machine Learning Repository. Available  
505 online: <http://archive.ics.uci.edu/ml>.
- 506 [48] G. Tzortzis, A. Likas, and G. Tzortzis, “The MinMax K-means clustering algorithm,”  
507 *Pattern Recognition*, vol. 47, no. 7, pp. 2505–2516, 2014.
- 508 [49] M. Halkidi, Y. Batistakis and M. Vazirgiannis, “On clustering validation tech-  
509 niques,” *Journal of Intelligent Information Systems Integrating Artificial Intelligence*  
510 *& Database Technologies*, vol. 17, no. 2, pp. 107–145, 2001.
- 511 [50] J. Demšar, “Statistical comparisons of classifiers over multiple data sets,” *Journal of*  
512 *Machine Learning Research*, vol. 7, no. 1, pp. 1–30, 2006.
- 513 [51] R. W. Johnson, “An introduction to the bootstrap,” *Teaching Statistics*, vol. 23, no. 2,  
514 pp. 49–54, 2001.
- 515 [52] Q. He, T. Shang, F. Zhuang, and Z. Shi, “Parallel extreme learning machine for  
516 regression based on mapreduce,” *Neurocomputing*, vol. 102, pp. 52–58, 2013.

# Figures

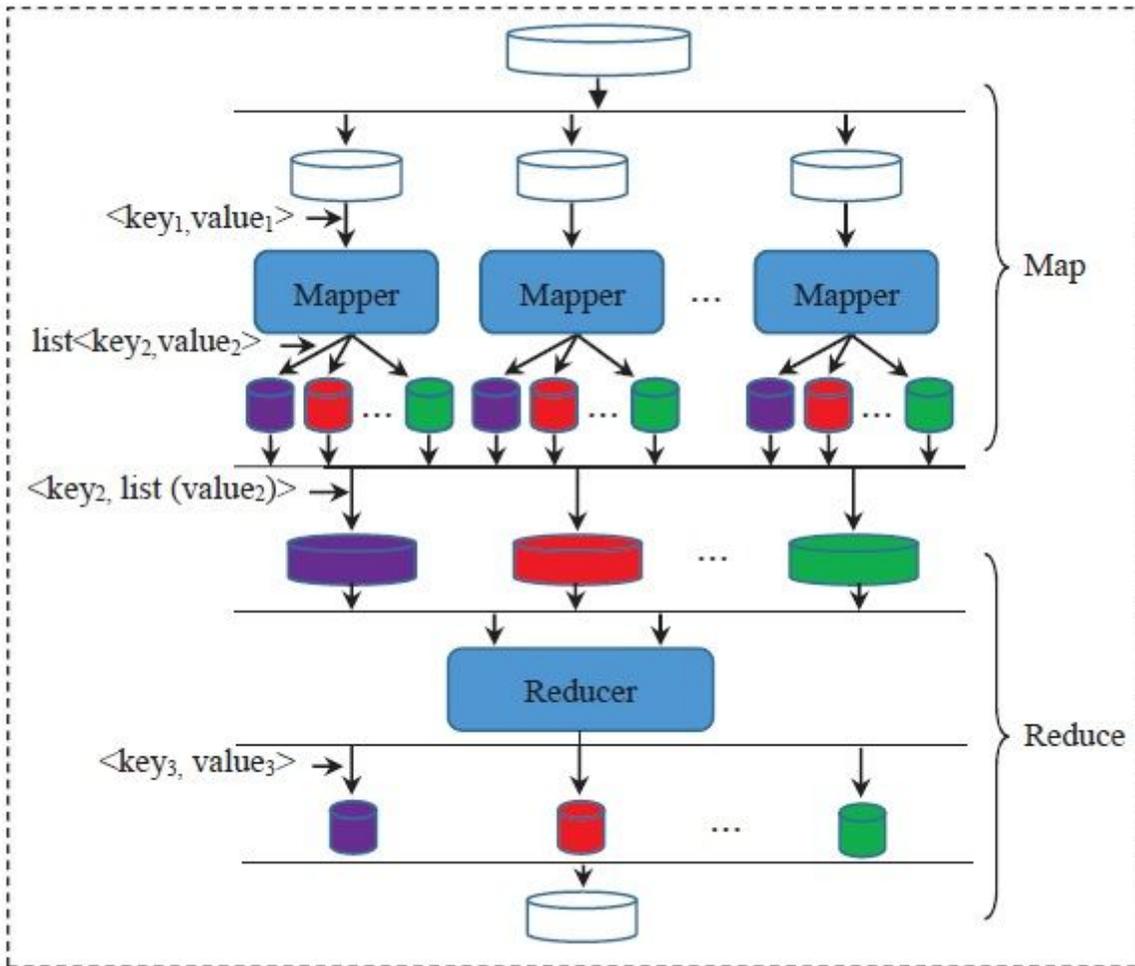
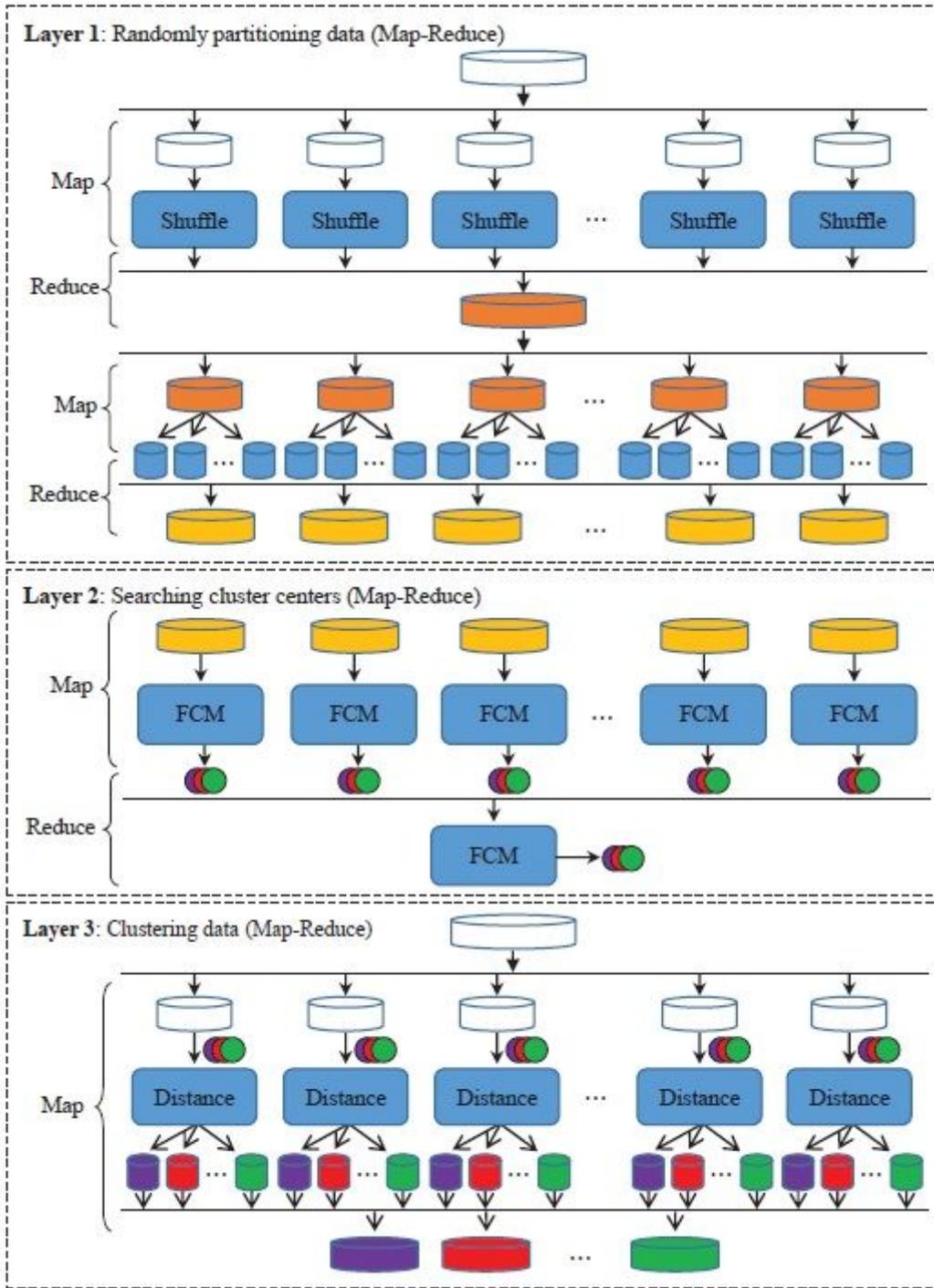


Figure 1

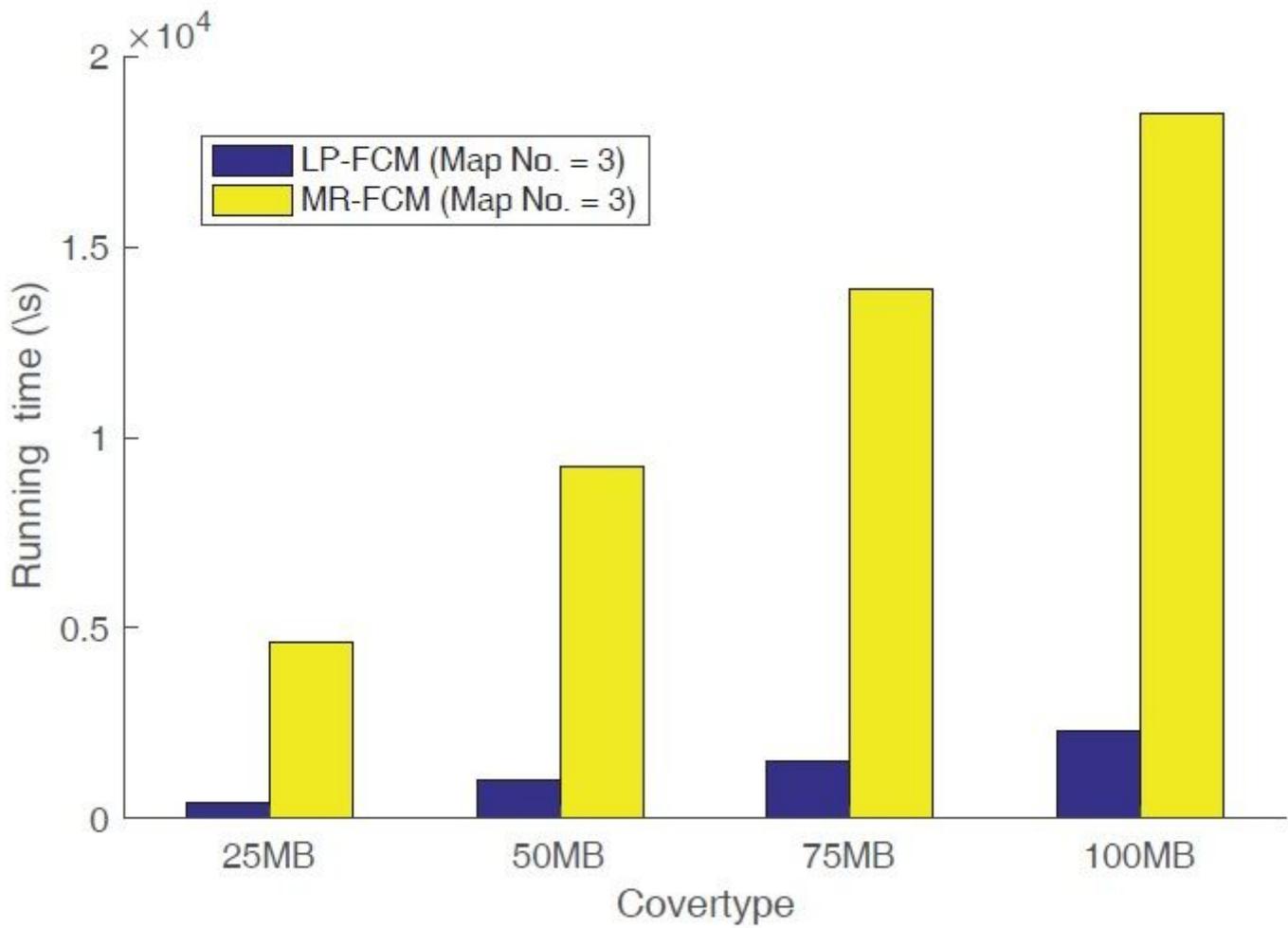
The detailed workflow of Map-Reduce.



**Figure 2**

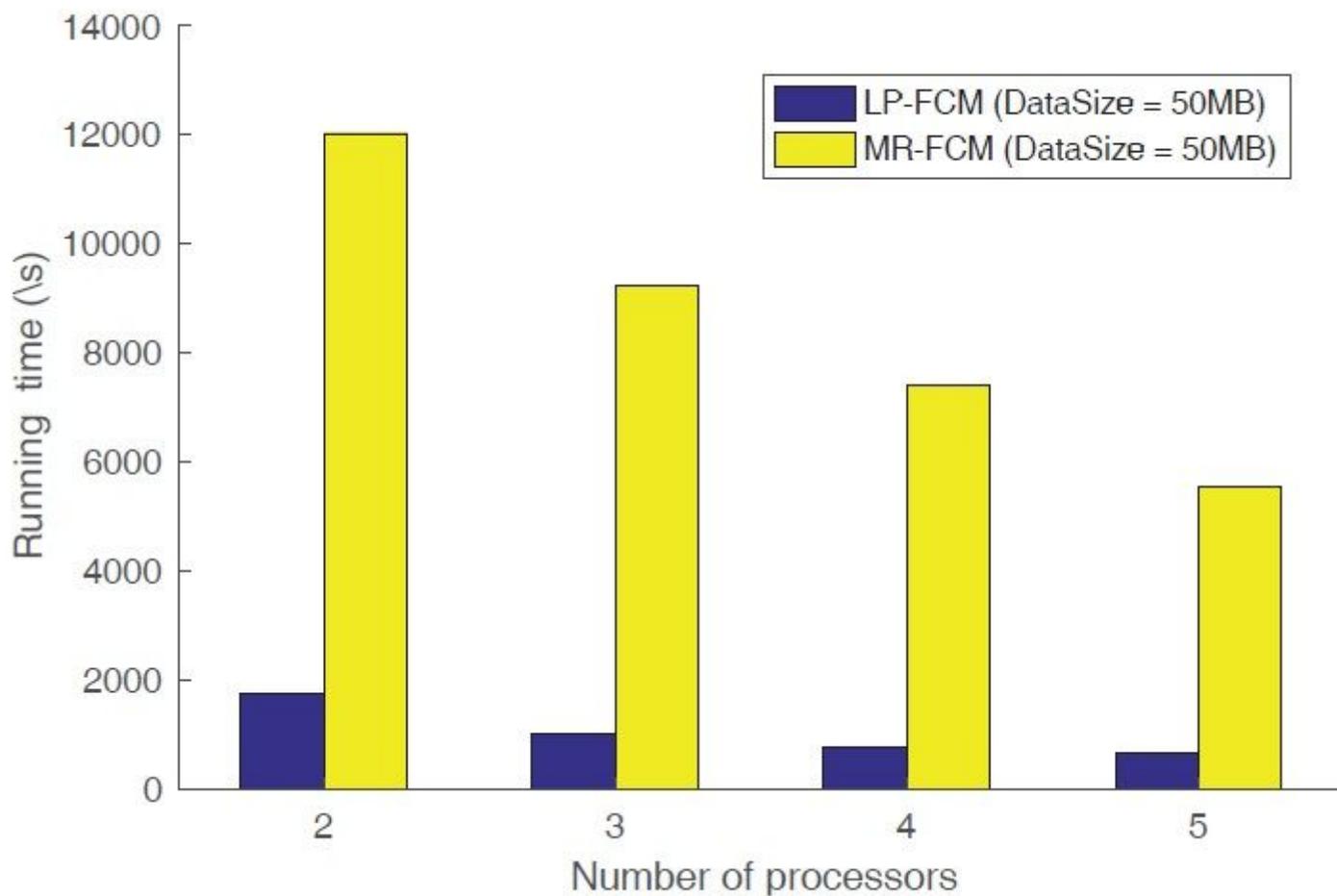
The main structure of LP-FCM algorithm. (1) Layer 1 contains two Map-Reduce jobs. The first job divides the original data (large hollow cylinder) into several sub-datasets (medium hollow cylinders) and shuffles each sub-dataset. Finally, the randomized data (large light red cylinder) are formed. The second job divides the randomized data into several sub-datasets (medium light red cylinders) and split each sub-dataset into several small data sets (small blue cylinders). Finally, the small data sets with same key are combined into one sub-dataset (medium yellow cylinder); (2) Layer 2 applies FCM to each sub-dataset (medium yellow cylinder) to confirm the centers (green purple/red/green small circle). The centers are

merged into a new dataset for applying FCM to determine the final centers. (3) Layer 3 labels each sample with the nearest center through calculating the distance between the sample and each center. Finally, the samples (small purple/red/green cylinder) belonging to the same category are gathered to a cluster (medium purple/red/green cylinder).



**Figure 3**

The execution time on different datasets.



**Figure 4**

The execution time on different processors.

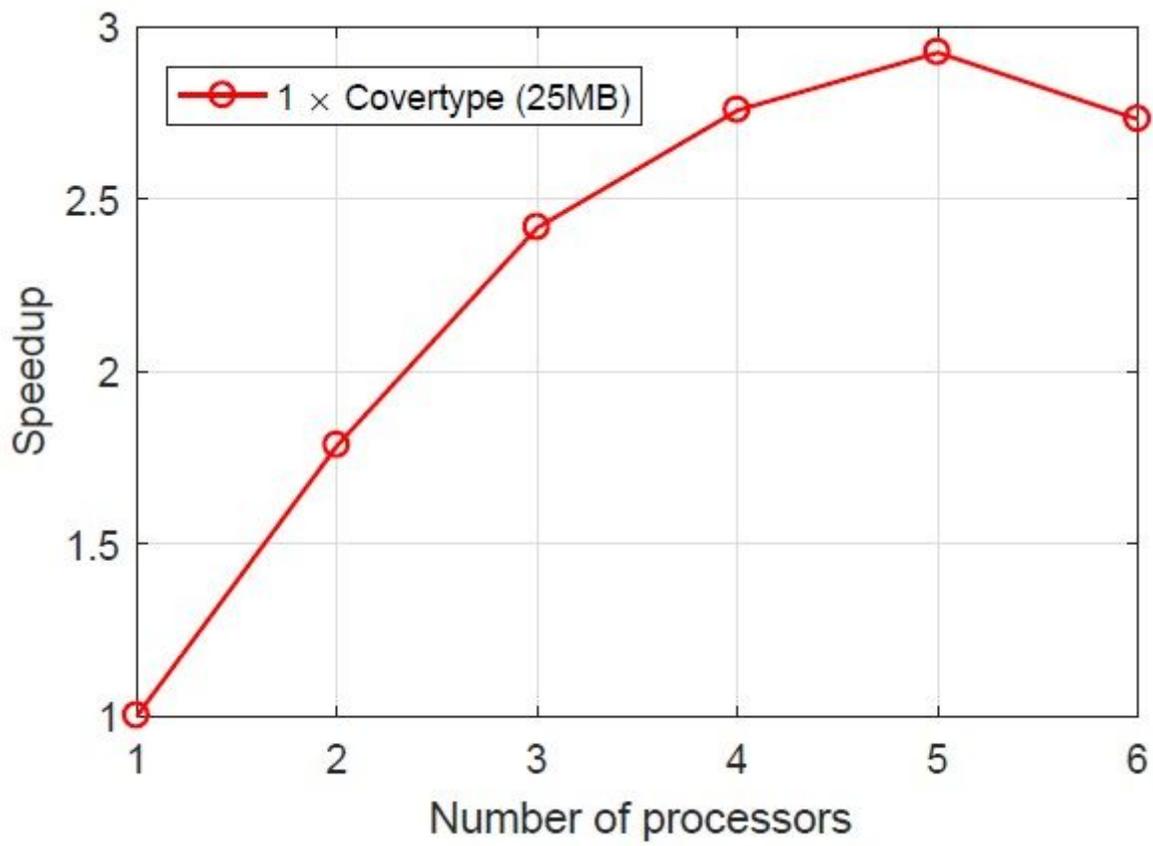


Figure 5

Speedup of the proposed algorithm for real datasets.

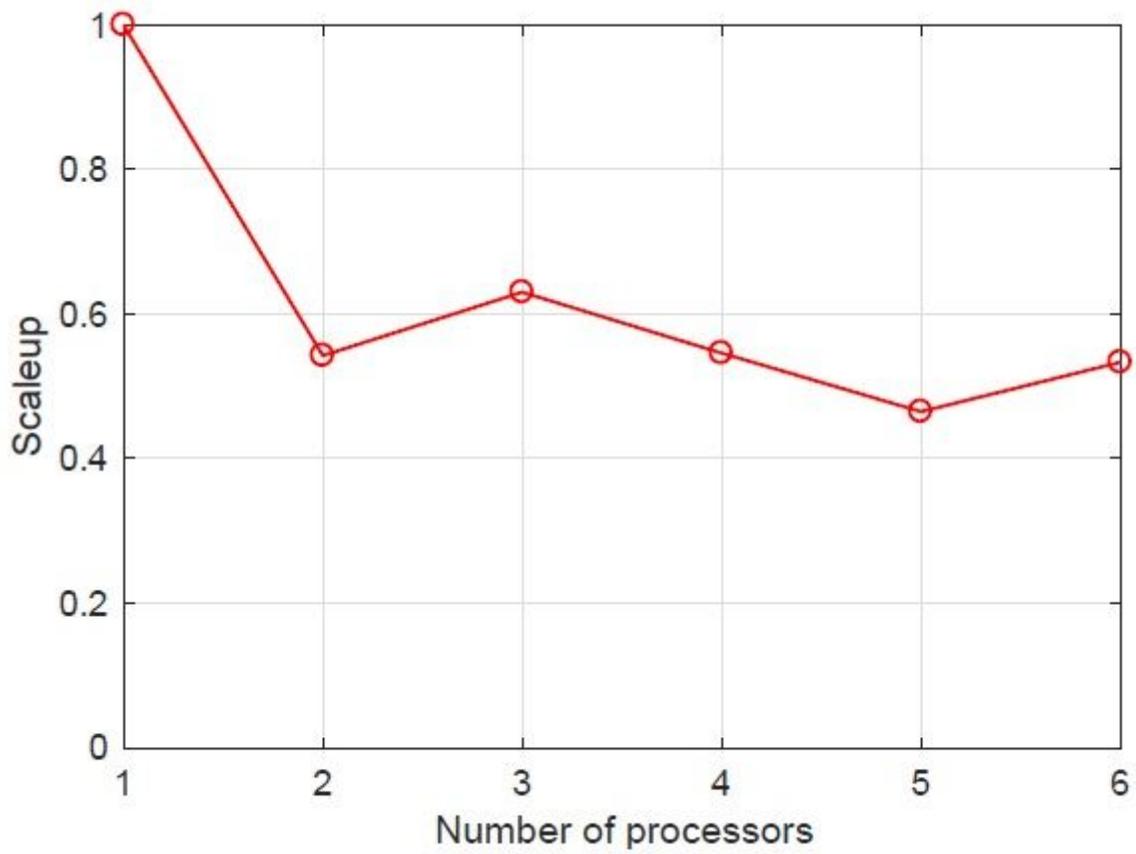


Figure 6

Scaleup of the proposed algorithm for real datasets.

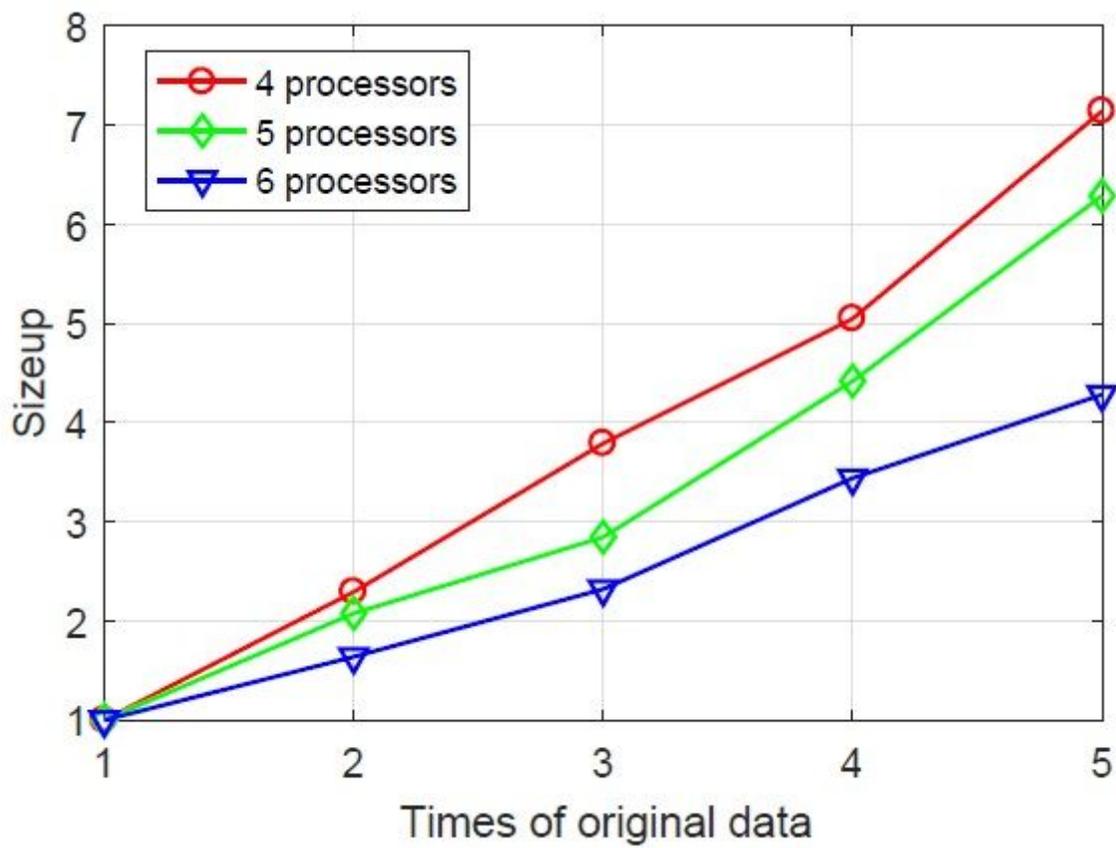


Figure 7

Sizeup of the proposed algorithm for real datasets.