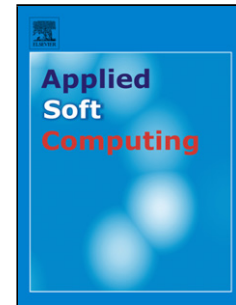


## Accepted Manuscript

Title: *k*-mw-modes: an algorithm for clustering categorical matrix-object data

Author: Fuyuan Cao Liqin Yu Joshua Zhexue Huang Jiye Liang



PII: S1568-4946(17)30191-6  
DOI: <http://dx.doi.org/doi:10.1016/j.asoc.2017.04.019>  
Reference: ASOC 4152

To appear in: *Applied Soft Computing*

Received date: 13-8-2016  
Revised date: 4-4-2017  
Accepted date: 11-4-2017

Please cite this article as: Fuyuan Cao, Liqin Yu, Joshua Zhexue Huang, Jiye Liang, *k*-mw-modes: an algorithm for clustering categorical matrix-object data, *Applied Soft Computing Journal* (2017), <http://dx.doi.org/10.1016/j.asoc.2017.04.019>

This is a PDF file of an unedited manuscript that has been accepted for publication. As a service to our customers we are providing this early version of the manuscript. The manuscript will undergo copyediting, typesetting, and review of the resulting proof before it is published in its final form. Please note that during the production process errors may be discovered which could affect the content, and all legal disclaimers that apply to the journal pertain.

# $k$ -mw-modes: an algorithm for clustering categorical matrix-object data

Fuyuan Cao<sup>a</sup>, Liqin Yu<sup>a</sup>, Joshua Zhexue Huang<sup>b</sup>, Jiye Liang<sup>a,\*</sup>

<sup>a</sup>Key Laboratory of Computational Intelligence and Chinese Information Processing of Ministry of Education, School of Computer and Information Technology, Shanxi University, Taiyuan 030006, China

<sup>b</sup>College of Computer Sciences & Software Engineering, Shenzhen University, Shenzhen 518060, China

---

## Abstract

In data mining, the input of most algorithms is a set of  $n$  objects and each object is described by a feature vector. However, in many real database applications, an object is described by more than one feature vector. In this paper, we call an object described by more than one feature vector as a matrix-object and a data set consisting of matrix-objects as a matrix-object data set. We propose a  $k$ -multi-weighted-modes (abbr.  $k$ -mw-modes) algorithm for clustering categorical matrix-object data. In this algorithm, we define the distance between two categorical matrix-objects and a multi-weighted-modes representation of cluster prototypes is proposed. We give a heuristic method to choose the locally optimal multi-weighted-modes in the iteration of the  $k$ -mw-modes algorithm. We validated the effectiveness and benefits of the  $k$ -mw-modes algorithm on the five real data sets from different applications.

*Keywords:* Categorical data, matrix-object,  $k$ -mw-modes algorithm

---

## 1. Introduction

In data mining, the input of an algorithm in most cases is a data set  $X$ , also called a table or matrix. The data set consists of  $n$  objects  $\{x_1, x_2, \dots, x_n\}$

---

\*Corresponding author

Email addresses: cfy@sxu.edu.cn (Fuyuan Cao), 1367545521@qq.com (Liqin Yu), zx.huang@szu.edu.cn (Joshua Zhexue Huang), 1jy@sxu.edu.cn (Jiye Liang)

and each object is described by  $m$  attributes  $\{A_1, A_2, \dots, A_m\}$  [1]. Most importantly, each object in  $X$  only corresponds with a feature vector  $(x_{i1}; x_{i2}; \dots; x_{im})$ ,  $i \in \{1, 2, \dots, n\}$ . However, in many real applications, a database often contains multiple tables. There are one-one, one-many or many-many relationships between two tables. Thus, an object usually corresponds with more than one transactional record. A real database application example from www.taobao.com is described in Table 1.

Table 1: A real application example from Taobao.

<i>User_ID</i>	<i>User_Sex</i>	<i>Age</i>	<i>Brand_ID</i>	<i>Brand_Name</i>	<i>Visit_Date</i>	<i>Visited_Times</i>
10944750	female	24	13451	WETHERM	06-04	8
			21110	SEMIR	06-07	1
			25687	JOSINY	05-08	11
			25687	JOSINY	05-15	4
			25687	JOSINY	05-17	2
			25687	JOSINY	06-06	2
			25687	JOSINY	06-15	2
			25687	JOSINY	07-02	5
			25687	JOSINY	07-25	2
			25687	JOSINY	08-09	3
			25687	JOSINY	08-13	1
8149250	male	29	18805	YEARCON	05-25	1
			18805	YEARCON	05-26	4
			18805	YEARCON	06-12	1
			21110	SEMIR	06-30	1
1832000	female	40	25687	JOSINY	08-14	18
			25687	JOSINY	08-15	3
			25687	JOSINY	06-26	5
			25687	JOSINY	08-11	1
			10151	OBERORA	08-03	1
			18805	YEARCON	07-26	5

10

There are two parts in Table 1. The left half describes the basic information

of users and the right one records that each user visited different brands in different time points, where the attribute *Visited\_Times* represents the visiting-times of a user on the same day for one brand. We call the left part as a master table and the right one as a detail table in database. Therefore, two parts in Table 1 exist a typical one-many relationship. Data in Table 1 have the following characteristics.

- **Correlation:** Data from the master table and the detail table maybe have some correlations. Users with different sex or age maybe have different preferences. For example, the female user of 24 years old from Table 1 visited the commodities that are usually used by most female users, such as *JOSINY* and *WETHERM*. However, the female user of 40 years old visited the commodities used by men or women, maybe because she needs to take after their families.
- **One-many:** Each user in the master table corresponds with more than one record in the detail table. Moreover, the number of brands visited by different users is often different in Table 1. For example, the user 10944750 has 11 records while the user 8149250 has 4 records.
- **Mixed:** In most cases, an object is described by categorical and numerical attributes together. For example, in the detail table, *Brand\_Name* is a categorical attribute while *Visited\_Times* is a numerical attribute.
- **Evolution:** Some attribute values will change as time goes on. For example, a user visits one brand repeatedly in this month, but the brand may be not visited by him or her in the next month. In other words, the change of a user's behaviour is a dynamic evolution process with time.

From the detail table, we can see clearly that every user visited one brand at least and a brand may be browsed by many users. Besides, a brand may be visited several times by a user in a day. Of course, it may also be visited many times by a user in several days. Obviously, if a user visited many times about a brand, he or she may be interested in this commodity. For example, for the

user 10944750, the *JOSINY* is visited in continuous four months, and there are several visiting times in every month. So, we can predict the user is likely very fond of the *JOSINY*. However, the *SEMIR* is visited only once by the user 10944750 in this data set, by which we know that the user may have less  
45 like about it compared with the *JOSINY*. Such a data representation shown in Table 1 is widespread in banking, insurance, telecommunication, retails, and medical databases. Therefore, it is necessary to develop a method that can discover user groups with different behaviour patterns from the detail table instead of the master table. Because the behaviour analysis can help managers  
50 obtain more valuable information for decision making.

Clustering is a widely used method to find different user groups in real applications [2] and the master table tends to be taken as its input. But the information in the master table cannot enough reflect the behaviour characteristics of a user. More importantly, in traditional clustering algorithms, the dissimilar-  
55 ity measure between two objects is based on the value difference of two feature vectors. For the detail table, each user has more than one transactional record. In other words, each user is described by multiple feature vectors. Therefore, some classical dissimilarity measures, such as Euclidean distance, Manhattan distance and Hamming distance, cannot be used to process this kind of data  
60 directly.

In the detail table, each user has multiple feature vectors, each of which is described by numerical and categorical attributes together in most cases. How to define a dissimilarity measure between two users is a very crucial problem, because it has direct effects on clustering results. For simplicity, in this paper, we  
65 only investigate the clustering algorithm for the detail table whose each record is described by categorical attributes. The  $k$ -modes algorithm [3] has realized the clustering of the categorical data sets compared with the  $k$ -means algorithm [4], but it still has some shortcomings. Only the data sets whose each object only contains one record can be clustered by the  $k$ -modes algorithm. Obviously,  
70 if the problem above wants to be solved with the  $k$ -modes algorithm, the data sets need to be compressed as the form that the algorithm required by selecting

an attribute value whose frequency is the highest. Thus, lots of information is at a loss in the data so that the clustering results are unfaithful.

Without loss of generality, a general description of detail information in Table 1 is illustrated as follows. Suppose that  $\mathbf{X} = \{X_1, X_2, \dots, X_n\}$  is a set of  $n$  objects described by  $m$  attributes  $\{A_1, A_2, \dots, A_m\}$ , where  $X_i = (X_{i1}; X_{i2}; \dots; X_{im})$  and  $X_{is} = [v_{i1s}, v_{i2s}, \dots, v_{ir_{is}}]'$ .

$r_i$  represents the number of records in  $X_i$  and  $v_{ijs}$  denotes the  $j$ th value of  $X_i$  on  $A_s$ . We call  $X_i$  as a matrix-object and  $\mathbf{X}$  as a matrix-object data set. Suppose that  $V^s$  represents the domain values of the attribute  $A_s$  in  $\mathbf{X}$  and  $V_{X_i}^{A_s}$  denotes a set of values on the attribute  $A_s$  for  $X_i$ . Obviously,  $\bigcup_{i=1}^n V_{X_i}^{A_s} = V^s$ . In traditional data representation, an object is only described by a feature vector or a record while a matrix-object is usually represented by multiple feature vectors or records. Therefore, a matrix-object is a general representation of a traditional object.

In this paper, we propose a new clustering algorithm, the  $k$ -mw-modes algorithm, to cluster categorical matrix-object data. The main contributions are summarized as follows:

- We define a new dissimilarity measure to calculate the distance between two categorical matrix-objects.
- We give a new representation and update way of the cluster centers to optimize the clustering process.
- We give a heuristic method to choose the cluster center of a set.
- We propose the  $k$ -mw-modes clustering algorithm to cluster categorical matrix-object data.
- Experimental results on the real data sets have shown the effectiveness of the  $k$ -mw-modes algorithm.

The rest of this paper is organized as follows. In Section 2, we propose the  $k$ -mw-modes algorithm. In section 3, we give a heuristic method to choose

100 the locally optimal multi-weighted-modes for the  $k$ -mw-modes algorithm. In Section 4, we show experimental results on the five real data sets from different applications. In Section 5, we review some related work. We give conclusions and future work in Section 6.

## 2. $k$ -multi-weighted-modes clustering

105 The  $k$ -modes clustering algorithm consists of three components: (1) representation of cluster centroids, (2) allocation of objects into clusters and (3) updates of cluster centroids. In this section, we present the  $k$ -mw-modes algorithm that uses the  $k$ -modes clustering process to cluster categorical matrix-object data. In this algorithm, we define a dissimilarity measure to calculate 110 the distance between two matrix-objects and give a kind of representation and update way of cluster centers.

### 2.1. Distance between two matrix-objects

Given two matrix-objects  $X_i$  and  $X_j$ , which are described by  $m$  categorical attributes  $\{A_1, A_2, \dots, A_m\}$ , the dissimilarity measure between  $X_i$  and  $X_j$  is 115 defined as

$$d(X_i, X_j) = \frac{1}{2} \sum_{s=1}^m \delta(X_{is}, X_{js}) \quad (1)$$

where

$$\delta(X_{is}, X_{js}) = \sum_{v \in V_{X_i}^{A_s} \cup V_{X_j}^{A_s}} \left| \frac{\sum_{p=1}^{r_i} f(v, v_{ips})}{r_i} - \frac{\sum_{q=1}^{r_j} f(v, v_{jq_s})}{r_j} \right| \quad (2)$$

and

$$f(x, y) = \begin{cases} 1, & \text{if } x == y. \\ 0, & \text{otherwise.} \end{cases} \quad (3)$$

Here,  $f(\cdot, \cdot)$  is a function whose value is 1 if two parameter values are equal, otherwise its value is 0.  $|\cdot|$  represents the absolute value of a value.

120 In addition, as  $0 \leq \delta(X_{is}, X_{js}) \leq 2$ , we add a normalization factor  $\frac{1}{2}$  in Eq.(1). We have  $\delta(X_{is}, X_{js}) = 2$  when  $V_{X_i}^{A_s} \cap V_{X_j}^{A_s} = \emptyset$ .

We can prove that the dissimilarity measure  $d(X_i, X_j)$  is a distance metric satisfying three properties as follows.

- 1) Nonnegativity:  $d(X_i, X_j) \geq 0$  and  $d(X_i, X_i) = 0$ ;
- 125 2) Symmetry:  $d(X_i, X_j) = d(X_j, X_i)$ ;
- 3) Triangle inequality:  $d(X_i, X_j) + d(X_j, X_k) \geq d(X_i, X_k)$ .

Obviously, we can easily prove the first two properties because  $|\cdot|$  is the symbol of an absolute value. The triangle inequality as the third property is verified as follows.

**Proof 1.** To prove the inequality  $d(X_i, X_j) + d(X_j, X_k) \geq d(X_i, X_k)$ , we only need to demonstrate

$$\delta(X_{is}, X_{js}) + \delta(X_{js}, X_{ks}) \geq \delta(X_{is}, X_{ks}), s \in \{1, 2, \dots, m\}.$$

With Eq.(2), the inequality above can be rewritten as

$$\begin{aligned} & \sum_{v \in V_{X_i}^{A_s} \cup V_{X_j}^{A_s}} \left| \frac{\sum_{p=1}^{r_i} f(v, v_{ips})}{r_i} - \frac{\sum_{q=1}^{r_j} f(v, v_{jq_s})}{r_j} \right| + \sum_{v \in V_{X_j}^{A_s} \cup V_{X_k}^{A_s}} \left| \frac{\sum_{q=1}^{r_j} f(v, v_{jq_s})}{r_j} - \frac{\sum_{l=1}^{r_k} f(v, v_{kls})}{r_k} \right| \\ = & \sum_{v \in V_{X_i}^{A_s} \cup V_{X_j}^{A_s} \cup V_{X_k}^{A_s}} \left| \frac{\sum_{p=1}^{r_i} f(v, v_{ips})}{r_i} - \frac{\sum_{q=1}^{r_j} f(v, v_{jq_s})}{r_j} \right| + \sum_{v \in V_{X_i}^{A_s} \cup V_{X_j}^{A_s} \cup V_{X_k}^{A_s}} \left| \frac{\sum_{q=1}^{r_j} f(v, v_{jq_s})}{r_j} - \frac{\sum_{l=1}^{r_k} f(v, v_{kls})}{r_k} \right| \\ = & \sum_{v \in V_{X_i}^{A_s} \cup V_{X_j}^{A_s} \cup V_{X_k}^{A_s}} \left( \left| \frac{\sum_{p=1}^{r_i} f(v, v_{ips})}{r_i} - \frac{\sum_{q=1}^{r_j} f(v, v_{jq_s})}{r_j} \right| + \left| \frac{\sum_{q=1}^{r_j} f(v, v_{jq_s})}{r_j} - \frac{\sum_{l=1}^{r_k} f(v, v_{kls})}{r_k} \right| \right) \\ \geq & \sum_{v \in V_{X_i}^{A_s} \cup V_{X_j}^{A_s} \cup V_{X_k}^{A_s}} \left| \frac{\sum_{p=1}^{r_i} f(v, v_{ips})}{r_i} - \frac{\sum_{q=1}^{r_j} f(v, v_{jq_s})}{r_j} + \frac{\sum_{q=1}^{r_j} f(v, v_{jq_s})}{r_j} - \frac{\sum_{l=1}^{r_k} f(v, v_{kls})}{r_k} \right| \\ = & \sum_{v \in V_{X_i}^{A_s} \cup V_{X_j}^{A_s} \cup V_{X_k}^{A_s}} \left| \frac{\sum_{p=1}^{r_i} f(v, v_{ips})}{r_i} - \frac{\sum_{l=1}^{r_k} f(v, v_{kls})}{r_k} \right| \\ = & \sum_{v \in V_{X_i}^{A_s} \cup V_{X_k}^{A_s}} \left| \frac{\sum_{p=1}^{r_i} f(v, v_{ips})}{r_i} - \frac{\sum_{l=1}^{r_k} f(v, v_{kls})}{r_k} \right|. \end{aligned}$$

130 The above proof verifies that the triangle inequality property holds on one attribute. It naturally can be generalized to multiple attributes. It follows that we have  $d(X_i, X_j) + d(X_j, X_k) \geq d(X_i, X_k)$ . Therefore, the dissimilarity measure  $d(\cdot, \cdot)$  is a distance metric.

Next, we give an example of calculating the distance between two matrix-objects as follows.

**Example 1.** Given two matrix-objects described by two categorical attributes whose values are represented by integers. The details are illustrated in Table 2.



Table 2: An example of two matrix-objects.

<i>ID</i>	<i>A</i> <sub>1</sub>	<i>A</i> <sub>2</sub>
<i>X</i> <sub>1</sub>	2	4
	2	5
	3	4
	4	4
	4	2
<i>X</i> <sub>2</sub>	2	4
	4	2
	5	2
	5	1

According to Eq.(1), the distance between  $X_1$  and  $X_2$  on the attribute  $A_1$  can be computed as  $|\frac{2}{5} - \frac{1}{4}| + |\frac{1}{5} - \frac{0}{4}| + |\frac{2}{5} - \frac{1}{4}| + |\frac{0}{5} - \frac{2}{4}| = 1$ . Similarly, the  
140 dissimilarity measure is  $\frac{11}{10}$  on the attribute  $A_2$ . Therefore we have  $d(X_1, X_2) = (1 + \frac{11}{10})/2 = \frac{21}{20}$ .

The dissimilarity measure  $d(\cdot, \cdot)$  is also a generalization of the simple matching dissimilarity measure that is used in the  $k$ -modes algorithm. In other words, if two matrix-objects have only one record respectively, their distance can be  
145 calculated by Eq.(1) as well. For instance, suppose that there are two objects  $X_1, X_2$  described by four attributes  $A_1, A_2, A_3, A_4$ ,  $X_1 = (a; c; d; c)$ ,  $X_2 = (a; d; d; b)$ . The distance between them by Eq.(1) is  $(|1 - 1| + (|1 - 0| + |0 - 1|) + |1 - 1| + (|1 - 0| + |0 - 1|))/2 = 2$ , and in the simple matching dissimilarity measure, it is  $0 + 1 + 0 + 1 = 2$ . So, the dissimilarity measure in the  $k$ -modes  
150 algorithm is a special case of  $d(\cdot, \cdot)$ .

With Eq.(1), we design an algorithm to calculate the distance between two matrix-objects. The details are described in *Algorithm 1*, which is named *ACDM* (an Algorithm of Calculating Distance between Matrix-objects).

## 2.2. Multi-weighted-modes as cluster centers

155 Suppose that  $\mathbf{X}$  is a matrix-object data set described by  $m$  categorical attributes, the cluster center of  $\mathbf{X}$  is defined as follows.

---

**Algorithm 1** The *ACDM*.

---

1: **Input:** -  $X_i, X_j$  : two matrix-objects described by  $m$  attributes;  
2: **Output:** - *distance*: the distance between  $X_i$  and  $X_j$ ;  
3: **Method:**  
4: *distance* = 0;  
5: **for**  $s = 1$  to  $m$  **do**  
6:   Obtain  $X_{is} = [v_{i1s}, v_{i2s}, \dots, v_{ir_i s}]'$  and  $X_{js} = [v_{j1s}, v_{j2s}, \dots, v_{jr_j s}]'$ ;  
7:    $V = V_{X_i}^{A_s} \cup V_{X_j}^{A_s} = \{v_1^s, v_2^s, \dots, v_u^s\}$ ;  
8:   *sum* = 0;  
9:   **for**  $t = 1$  to  $u$  **do**  
10:      $s1 = 0, s2 = 0$ ;  
11:     **for**  $p = 1$  to  $r_i$  **do**  
12:       **if**  $v_t^s == v_{ips}$  **then**  
13:          $s1 = s1 + 1$ ;  
14:       **end if**  
15:     **end for**  
16:     **for**  $q = 1$  to  $r_j$  **do**  
17:       **if**  $v_t^s == v_{jq s}$  **then**  
18:          $s2 = s2 + 1$ ;  
19:       **end if**  
20:     **end for**  
21:     *sum* = *sum* +  $|s1/r_i - s2/r_j|$ ;  
22:   **end for**  
23:   *distance* = *distance* + *sum*/2;  
24: **end for**  
25: return *distance*;

---

**Definition 1.** Let  $V_{X_i}^{A_s} = \{v_1^s, v_2^s, \dots, v_{u'_s}^s\}$  be the domain values of  $X_i$  on the attribute  $A_s$ , the frequency of  $v_u^s$  ( $1 \leq u \leq u'_s$ ) in  $X_i$  is defined as

$$f'_i(v_u^s) = \sum_{p=1}^{r_i} f(v_u^s, v_{ips}), \quad (4)$$

160 With Eq.(4), we can easily obtain the frequency of each attribute value in  $X_i$  on  $A_s$ . Thus,  $X_{is}$  can be represented as  $X_{is} = \{(v_u^s, f'_i(v_u^s)) | v_u^s \in V_{X_i}^{A_s}\}$  or  $X_{is} = \{(v_u^s, \frac{f'_i(v_u^s)}{r_i}) | v_u^s \in V_{X_i}^{A_s}\}$ . We call  $\frac{f'_i(v_u^s)}{r_i}$  as the weight of  $v_u^s$  in  $X_i$ .

**Definition 2.** Let  $\mathbf{X} = \{X_1, X_2, \dots, X_n\}$  be a set containing  $n$  matrix-objects and  $Q = (Q_{A_1}; Q_{A_2}; \dots; Q_{A_m})$  be a matrix-object. They are all described by  
165  $m$  categorical attributes.  $Q$  is the multi-weight-modes or the center of  $\mathbf{X}$  if  $Q$  minimizes the following function

$$F(\mathbf{X}, Q) = \sum_{i=1}^n d(X_i, Q), \quad (5)$$

where  $X_i \in \mathbf{X}$  and  $d(X_i, Q)$  can be calculated by Eq.(1).

To minimize  $F(\mathbf{X}, Q)$ , we only need to minimize  $\sum_{i=1}^n \delta(X_{is}, Q_{A_s})$ , the sum of the distance between  $Q$  and each matrix-objects in  $\mathbf{X}$  on the attribute  $A_s$   
170 where  $s \in \{1, 2, \dots, m\}$ . As the attribute values in  $Q_{A_s}$  must be from the values in  $V^s$  that represents the domain values of  $\mathbf{X}$  on  $A_s$ , the number of categorical values in  $Q_{A_s}$  is between 1 and  $|V^s|$ . According to Definition 1, we can compute  $\sum_{i=1}^n (f'_i(v)) (v \in V^s)$  in  $\mathbf{X}$ . If we choose  $u_s$  values  $\{v_1^s, v_2^s, \dots, v_{u_s}^s\}$  from  $V^s$  as the values of  $Q_{A_s}$ , there are  $C_{|V^s|}^{u_s}$  combinations. For a given combination,  
175 its component is represented by  $(v_j^s, \sum_{i=1}^n (f'_i(v_j^s)))$  where  $j \in \{1, 2, \dots, u_s\}$ . It follows that the total number of possible sets for  $Q_{A_s}$  is  $\sum_{u_s=1}^{|V^s|} C_{|V^s|}^{u_s}$ . Therefore, we need to traverse every combination to find a  $Q_{A_s}$ , which minimizes  $\sum_{i=1}^n \delta(X_{is}, Q_{A_s})$ . A global optimization algorithm for finding multi-weighted-modes is described in *Algorithm 2*, which is named *GAFMWM* (A Global Algorithm of Finding Multi-Weighted-modes).  
180

### 2.3. The $k$ -mw-modes algorithm

Given Eq.(1) as the distance measure between two matrix-objects, the  $k$ -mw-modes algorithm for clustering a matrix-object set  $\mathbf{X} = \{X_1, X_2, \dots, X_n\}$

---

**Algorithm 2** The *GAFMWM*.

---

1: **Input:** -  $\mathbf{X}$ : A set of  $n$  matrix-objects described by  $m$  attributes;  
2: **Output:** -  $Q$  : The cluster center of  $\mathbf{X}$ ;  
3: **Method:**  
4:  $Q = \emptyset$ ;  
5: **for**  $s = 1$  to  $m$  **do**  
6:   Generate a set  $q^s = \{q_s^1, q_s^2, \dots, q_s^{2^{|V^s|}-1}\}$  in  $V^s$  by binomial theorem;  
7:   **for**  $j = 1$  to  $2^{|V^s|} - 1$  **do**  
8:      $Q_{A_s}^j = \emptyset$ ;  
9:     **for**  $p = 1$  to  $|q_s^j|$  **do**  
10:       $Q_{A_s}^j = Q_{A_s}^j \cup \{(v_p^s, \sum_{i=1}^n f'_i(v_p^s)) | v_p^s \in q_s^j\}$ ;  
11:     **end for**  
12:     Computer  $F_j = \sum_{i=1}^n \delta(X_{is}, Q_{A_s}^j)$  by Eq.(2);  
13:     **if**  $j = 1$  **then**  
14:       set  $\text{minValue} = F_j$ ,  $Q_{A_s} = Q_{A_s}^j$ ;  
15:     **else if**  $F_j < \text{minValue}$  **then**  
16:       set  $\text{minValue} = F_j$ ,  $Q_{A_s} = Q_{A_s}^j$ ;  
17:     **end if**  
18:   **end for**  
19:   set  $Q = Q \cup Q_{A_s}$ ;  
20: **end for**  
21: return  $Q$ ;

---

into  $k (\ll n)$  clusters minimizes the following objective function

$$F'(W, Q) = \sum_{l=1}^k \sum_{i=1}^n \omega_{li} d(X_i, Q_l),$$

subject to

$$\omega_{li} \in \{0, 1\}, 1 \leq l \leq k, 1 \leq i \leq n, \quad (6)$$

$$\sum_{l=1}^k \omega_{li} = 1, 1 \leq i \leq n, \quad (7)$$

$$0 < \sum_{i=1}^n \omega_{li} < n, 1 \leq l \leq k, \quad (8)$$

185 where  $W = [\omega_{li}]$  is a  $k$ -by- $n$   $\{0, 1\}$  matrix in which  $\omega_{li} = 1$  indicates that object  $X_i$  is allocated to cluster  $l$ ,  $Q = [Q_1, Q_2, \dots, Q_k]$  in which  $Q_l \in Q$  is the multi-weighted modes of cluster  $l$ .

$F'(W, Q)$  can be solved with an iterative process to solve two subproblems iteratively until the process converges. The first step is to fix  $Q = Q^t$  at iteration  $t$  and solve the reduced problem  $F'(W, Q^t)$  with Eq.(1) to find  $W^t$  that minimizes  $F'(W, Q^t)$ . The second step is to fix  $W^t$  and solve the reduced problem  $F'(W^t, Q)$  by using algorithm *GAFMWM* to find  $Q^{t+1}$  that minimizes  $F'(W^t, Q)$ . We give the description of the algorithm in *Algorithm 3*, which is named as the  $k$ -mw-modes algorithm.

195 The computation complexity of the  $k$ -mw-modes algorithm is analyzed as follows.

- The computation complexity for calculation of the distance between two matrix-objects on  $A_s$  is  $\mathcal{O}(|V^s|)$ . The computation complexity of the distance between two matrix-objects in  $m$  attributes is  $\mathcal{O}(m \times |V'|)$ , where  
200  $|V'| = \max\{|V^s|, 1 \leq s \leq m\}$ .
- Updating cluster centers. The main goal of updating cluster centers is to find the multi-weighted-modes in each cluster according to the partition matrix  $W$ . The computational complexity for this step is  $\mathcal{O}(km \times 2^{|V'|})$ , where  $|V'| = \max\{|V^s|, 1 \leq s \leq m\}$ .

205 If the clustering process needs  $t$  iterations to converge, the total computational complexity of the  $k$ -mw-modes algorithm is  $\mathcal{O}(nmtk \times 2^{|V'|})$ , where  $|V'| = \max\{|V^s|, 1 \leq s \leq m\}$ . It is obviously that the time complexity of the proposed algorithm increases linearly as the number of objects and clusters increases, and increases exponentially with the increasing of the number  
210 of attribute values. The space complexity of the  $k$ -mw-modes algorithm is  $\mathcal{O}((n+k) \sum_{s=1}^m |V^s|)$ .

---

**Algorithm 3** The  $k$ -mw-modes Algorithm.

---

1: **Input:**

2: -  $\mathbf{X}$ : a data set of  $n$  matrix-objects described by  $m$  attributes;

3: -  $k$ : the number of clusters need to be clustered;

4: -  $\varepsilon$ : A threshold;

5: -  $idCenters$ : the  $id$  of  $k$  initial centers;

6: **Output:**

7: -  $cid$ : the labels of all objects after clustering;

8: -  $num$ : the iterations;

9: **Method:**

10: Let  $Q$  store the  $k$  initial centers by the indexes in  $idCenters$ ;

11: value=0, num=0;

12: **while** num  $\leq$  100 **do**

13:   newvalue=0;

14:   **for**  $i = 1$  to  $n$  **do**

15:     **for**  $j = 1$  to  $k$  **do**

16:       Calculate the distance between  $i$ th object and  $j$ th clustering center  
by Eq.(1);

17:     **end for**

18:     Arrange the  $i$ th object to the  $l$ th cluster if  $l = \arg \min_{j=1}^k \{d(X_i, Q_j)\}$ .

19:     newvalue=newvalue+ $\min_{j=1}^k \{d(X_i, Q_j)\}$ ;

20:   **end for**

21:   If  $|newvalue - value| \leq \varepsilon$ , *break*; Else  $value = newvalue$  and  $num = num + 1$ ;

22:   **for**  $i = 1$  to  $k$  **do**

23:     Update the cluster centers  $Q$  with *Algorithm 2*;

24:   **end for**

25: **end while**

---

**Theorem 1.** *The  $k$ -mw-modes algorithm converges to a local minimal solution in a finite number of iterations.*

**Proof 2.** We note that the number of possible values for the center of a cluster  
 215 is  $N = \prod_{s=1}^m \sum_{u_s=1}^{|V^s|} C_{|V^s|}^{u_s}$ , where  $|V^s|$  is the number of unique values on  $A_s$  and  
 $C_{|V^s|}^{u_s}$  is the number of combinations of choosing  $u_s$  values from a set of  $|V^s|$   
 values. To divide a data set into  $k$  clusters, the number of possible partitions  
 is finite. We can show that each possible partition only occurs once in the  
 clustering process. Let  $W^h$  be a partition at iteration  $h$ . We can obtain  $Q^h$  that  
 220 depends on  $W^h$ .

Suppose that  $W^{h_1} = W^{h_2}$ , where  $h_1$  and  $h_2$  are two different iterations, i.e.,  
 $h_1 \neq h_2$ . If  $Q^{h_1}$  and  $Q^{h_2}$  are obtained from  $W^{h_1}$  and  $W^{h_2}$ , respectively, then  
 $Q^{h_1} = Q^{h_2}$  since  $W^{h_1} = W^{h_2}$ . Therefore, we have

$$F'(W^{h_1}, Q^{h_1}) = F'(W^{h_2}, Q^{h_2}).$$

However, the value of the objective function  $F'(\cdot, \cdot)$  generated by the  $k$ -MW-  
 modes algorithm is strictly decreasing.  $h_1$  and  $h_2$  must be two consecutive it-  
 erations in which the clustering result is no longer change and the clustering  
 process converges. Therefore, the  $k$ -mw-modes algorithm converges in a finite  
 225 number of iterations.

### 3. A heuristic method for updating cluster centers

The *GAFMWM* for finding cluster centers is not efficient if the number  
 of domain values is very large. In this section, we give a heuristic method of  
 updating cluster centers in the  $k$ -mw-modes clustering process. For  $X_i, X_j \in \mathbf{X}$ ,  
 230 we have  $V_{X_i}^{A_s} = V_{X_j}^{A_s}$  or  $V_{X_i}^{A_s} \neq V_{X_j}^{A_s}$  on the attribute  $A_s$ . Even if  $V_{X_i}^{A_s} = V_{X_j}^{A_s}$ , the  
 frequency of the same attribute value may be different in  $X_i$  and  $X_j$ , because a  
 value maybe appears more than once in a given matrix-object. The higher the  
 frequency of a value in a given matrix-object is, the more the possibility of the  
 value as the cluster center is. In order to describe the possibility of a value as  
 235 the cluster center, the weight of a value is defined as follows.

**Definition 3.** Let  $V^s = \{v_1^s, v_2^s, \dots, v_{u_s}^s\}$  be the domain values of  $\mathbf{X}$  on the  
 attribute  $A_s$ . For any  $u \in \{1, 2, \dots, u_s\}$ , the weight of  $v_u^s$  in  $\mathbf{X}$  is defined as

$$\omega(v_u^s) = \frac{1}{n} \sum_{i=1}^n \frac{f'_i(v_u^s)}{r_i}. \quad (9)$$

Obviously, we can obtain the weight of each value in  $V^s$  with Eq.(9) and  
 240 arrange them in the descending order of the weight. Suppose that  $V^s =$   
 $\{v_{q_1}^s, v_{q_2}^s, \dots, v_{q_{u'_s}}^s\}$  and  $\omega(v_{q_1}^s) \geq \omega(v_{q_2}^s) \geq \dots \geq \omega(v_{q_{u'_s}}^s)$ . According to Def-  
 inition 1, we have  $\sum_{i=1}^n f'_i(v_{q_1}^s), \sum_{i=1}^n f'_i(v_{q_2}^s), \dots, \sum_{i=1}^n f'_i(v_{q_{u'_s}}^s)$ . If there are  
 $u_s (1 \leq u_s \leq u'_s)$  values in  $Q_{A_s}$ ,  $Q_{A_s} = \{(v_j^s, \sum_{i=1}^n f'_i(v_j^s)) | q_1 \leq j \leq q_{u_s}\}$ . Here,  
 in order to increase efficiency, we set  $u_s = \text{round}(\frac{\sum_{i=1}^n |V_{X_i}^{A_s}|}{n})$  for  $Q_{A_s}$ . By the  
 245 way above, the centers on other attributes can be found as well. The heuristic  
 algorithm is described in *Algorithm 4*, which is called *HAFMWM* (A Heuristic  
 Algorithm of Finding Multi-Weighted-modes).

The time complexity of *HAFMWM* is  $\mathcal{O}(km|V'|)$ , where  $|V'| = \max\{|V^s|, 1 \leq$   
 $s \leq m\}$ . Obviously, the computation complexity of *HAFMWM* is less than  
 250 *GAFMWM*.

Next, we give an example to describe the process of finding cluster centers  
 by *HAFMWM* as follows.

**Example 2.** Suppose that  $\mathbf{X} = \{X_1, X_2, X_3, X_4\}$  is a set described by one cat-  
 egorical attribute  $A_1$ .  $X_1, X_2, X_3, X_4$  are matrix-objects and contain five, four,  
 four and two records, respectively. The details are described in Table 3.

Table 3: An example data set for finding center.

<i>ID</i>	$A_1$	<i>ID</i>	$A_1$	<i>ID</i>	$A_1$	<i>ID</i>	$A_1$
$X_1$	3	$X_2$	3	$X_3$	3	$X_4$	3
	3		6		5		7
	5		4		6		
	4		3		4		
	4						

255

From Table 3, we can obtain the domain values  $V^1 = \{3, 4, 5, 6, 7\}$  on the  
 attribute  $A_1$ . According to the Eqs.(3)-(4) and Eq.(9), the weight of each value



---

**Algorithm 4** The *HAFMWM*.

---

1: **Input:** -  $\mathbf{X}$  : A set of  $n$  matrix-objects described by  $m$  attributes;  
2: **Output:** -  $Q$  : A cluster center of  $\mathbf{X}$ ;  
3: **Method:**  
4: **for**  $s = 1$  to  $m$  **do**  
5:    $sum=0, Q_{A_s} = \emptyset$ ;  
6:   **for**  $i = 1$  to  $n$  **do**  
7:      $sum = sum + |V_{X_i}^{A_s}|$ ;  
8:   **end for**  
9:    $u_s = round(sum/n)$ ;  
10:   **for**  $i = 1$  to  $|V^s|$  **do**  
11:     Calculate the weight  $\omega(v_i^s)$  of  $v_i^s \in V^s$  by Eq.(9);  
12:   **end for**  
13:   Arrange  $\{v_1^s, v_2^s, \dots, v_{|V^s|}^s\}$  in the descending order of their weight;  
14:   **for**  $p = 1$  to  $u_s$  **do**  
15:     Compute  $\sum_{i=1}^n f'_i(v_p^s)$ ;  
16:      $Q_{A_s} = Q_{A_s} \cup \{(v_p^s, \sum_{i=1}^n f'_i(v_p^s))\}$ ;  
17:   **end for**  
18:    $Q = Q \cup Q_{A_s}$ ;  
19: **end for**  
20: return  $Q$ ;

---

can be calculated as follows:

$$\begin{aligned}\omega(3) &= \frac{1}{4} \left( \frac{f'_1(3)}{r_1} + \frac{f'_2(3)}{r_2} + \frac{f'_3(3)}{r_3} + \frac{f'_4(3)}{r_4} \right) \\ &= \frac{1}{4} \left( \frac{2}{5} + \frac{2}{4} + \frac{1}{4} + \frac{1}{2} \right) \\ &= \frac{1.65}{4},\end{aligned}$$

In the same way, we have  $\omega(4) = \frac{0.90}{4}$ ,  $\omega(5) = \frac{0.45}{4}$ ,  $\omega(6) = \frac{0.5}{4}$ ,  $\omega(7) = \frac{0.5}{4}$ . Clearly,  $\omega(3) > \omega(4) > \omega(6) = \omega(7) > \omega(5)$ . And  $u_1 = round(\frac{3+3+4+2}{4}) = 3$ , so we choose 3, 4, 6 as values of  $Q_{A_1}$ . Furthermore, we can compute  $\sum_{i=1}^4 f'_i(3) = 6$ ,  $\sum_{i=1}^4 f'_i(4) = 4$ ,  $\sum_{i=1}^4 f'_i(6) = 2$ . Therefore, the center of  $\mathbf{X}$  can be represented as  $((3, 6), (4, 4), (6, 2))$ .

260

#### 4. Experiments on real data

In this section, we mainly make some experiments on the five real data sets, Microsoft Web data, Market Basket data, Alibaba data, Musk data and Movie-lens data, to evaluate the effectiveness of the proposed algorithm. We firstly  
 265 describe the preprocessing process of the five data sets. Then five evaluation indexes are introduced. Finally, we show the comparison results of the  $k$ -mw-modes algorithm with other algorithms and discuss the impact of the parameter  $\varepsilon$  on the clustering performance.

##### 4.1. Data preprocessing

270 To our best knowledge, open matrix-object data with label information are very rare. To tackle this problem, we need to conduct data preprocessing for the given real data sets, because the data sets clustered by clustering algorithms are generally supposed to exist some structure or distribution.

To obtain the structure of a given matrix-object data set, we use the mul-  
 275 tidimensional scaling technique [5] to visualize the data. The main goal of the technique is to obtain a configuration of  $n$  points (rows) in  $P$  dimensions (cols) by passing the  $n$ -by- $n$  distance matrix obtained by Eq.(1) to the function *mdscale* from *MATLAB*. The Euclidean distances between  $n$  points approximate a monotonic transformation of the corresponding dissimilarities in the  
 280  $n$ -by- $n$  distance matrix. Therefore, we can visualize  $n$  points to reflect the distribution of the data. To visualize the data, we set  $P = 2$ .

In most cases, the distribution of a real data set is often disordered. By the visualization technology, we can delete some points to get the relative clear structure of the data. From the visual figure, we can intuitively find the number  
 285 of clusters and obtain the label information of every matrix-object. Thus, we can use external evaluation indexes to evaluate the clustering performance of the  $k$ -mw-modes algorithm. Below is the preprocessing process of the five real data sets.

#### 4.1.1. Microsoft Web data

Microsoft Web data set downloaded from UCI is created by sampling and  
 290 processing the www.microsoft.com logs. It records the visiting of the web sites in  
 one week timeframe in February 1998 by 32711 anonymous, randomly-selected  
 users. For each user, it is described by two attributes, User\_Id and Web\_Id, and  
 more than one web site are visited. So, each user is a matrix-object and the  
 295 data set can be used to evaluate the  $k$ -mw-modes algorithm.

The preprocessing process is as follows: firstly eliminate the objects that  
 visit less than seven web sites to generate a temporary set; secondly visualize  
 the temporary set in the coordinate system [5] and select the objects whose  
 abscissa values are in the position of  $x < -0.1$  or  $x > 0.1$  to form a new set of  
 300 2401 objects; finally visualize the new set shown in Fig.1.

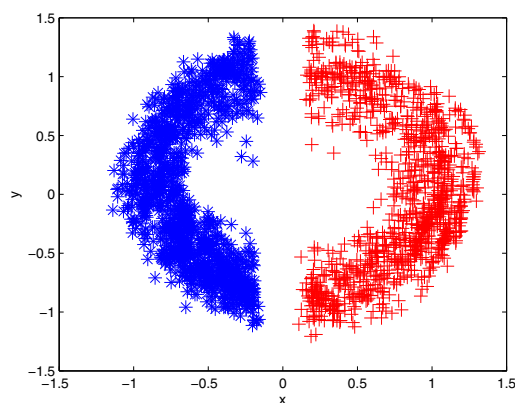


Figure 1: The distribution of Microsoft Web data after preprocessing.

Obviously, Microsoft Web data set can be divided into 2 clusters.

#### 4.1.2. Market Basket data

Market Basket data downloaded from Data website<sup>1</sup> record the transactions  
 of 1001 customers, each of which is described by four attributes, Customer\_Id,

<sup>1</sup><http://www.datatang.com/datares/go.aspx?dataid=613168>

305 Time, Product\_Name and Product\_Id. Here, we just need to take Customer\_Id  
 and Product\_Id into account and ignore the other attributes, because all cus-  
 tomers have the same values on the attribute Time and Product\_Name can be  
 replaced completely by Product\_Id. In addition, the prominent characteristic  
 for Market Basket data is that every customer in it has 7 transactional records.  
 310 Therefore, each customer is a typical matrix-object.

Following is the process of Market Basket data preprocessing. By visualizing  
 the data with multidimensional scaling technique, we select some objects who  
 locate the position of  $x < -0.2, y < 0.5$  or  $x > 0, y < -0.1$  or  $x > -0.3, y > 0.7$   
 or  $x > 0.3, y > 0.1$  in the coordinate system to form a new data set of 900  
 315 objects. The distribution of the new data set is shown in Fig.2.

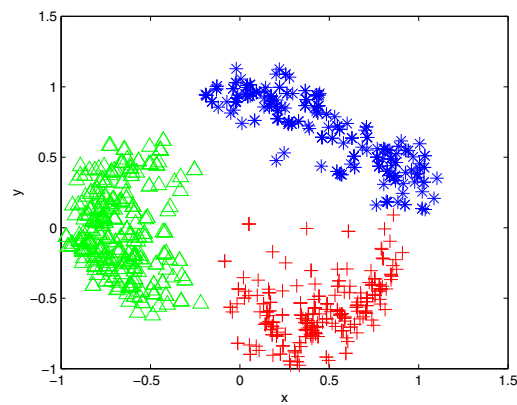


Figure 2: The distribution of Market Basket data after preprocessing.

It is clearly that Market Basket data can be divided into 3 clusters.

#### 4.1.3. Alibaba data

Alibaba data downloaded from the competition website<sup>2</sup> describe user's be-  
 haviour of visiting brands. It records 182880 visiting records of 884 users who  
 320 are described by four attributes, User\_Id, Time, Action\_type and Brand\_Id. In

<sup>2</sup><http://102.alibaba.com/competition/addDiscovery/index.htm>

this experiment, we only consider the attribute User.Id and Brand.Id. We can see that Alibaba data set is a matrix-object data set because every user visited more than one brand.

The preprocessing process of Alibaba data is described as follows: firstly  
 325 visualize the data in the coordinate system; secondly eliminate the objects whose abscissa values are in the range of  $-0.2 < x < 0.2$  or  $x > 0.2, 0 < y < 0.2$  to obtain a new set of 793 objects; finally visualize the new set in Fig.3.

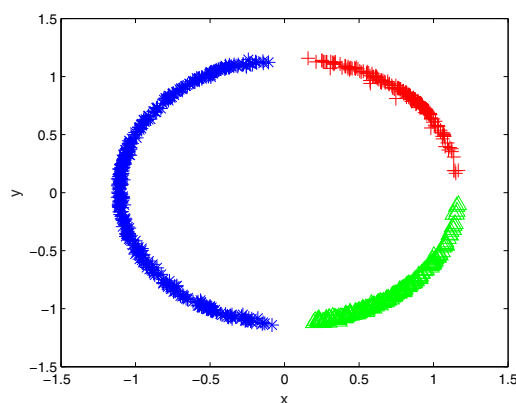


Figure 3: The distribution of Alibaba data after preprocessing.

Clearly, Alibaba data are divided into 3 clusters according to Fig.3.

#### 4.1.4. Musk data

330 Musk data describe a set of 92 objects, each of which represents a molecule and is described by 167 attributes. This data set can be downloaded from UCI [6] and aims at predicting whether new molecules will be musks or non-musks. The fact that 476 instances or records are contained in Musk data results in more than one record for an object. That is to say, Musk data set is a matrix-  
 335 object data set. Furthermore, it has been divided into 2 clusters by the human experts that the 47 molecules are judged to be musks while the remaining 45 molecules are non-musks. In this experiment, we consider attribute values on Musk data as categorical values. Therefore, we can cluster Musk data directly

without preprocessing.

#### 340 4.1.5. MovieLens data

MovieLens data are collected and made from the MovieLens website<sup>3</sup> and they contain three different size parts, MovieLens 100k, MovieLens 1M and MovieLens 10M. In this experiment, we make some experiments on the MovieLens 1M data, which contain three files including movies data, ratings data and  
345 users data. Movies data and users data are not considered because they only describe the fundamental information of movies and users respectively.

In this experiment we only employ the ratings data set which records 1000209 ratings of 3900 movies made by 6040 users selected randomly. And it is described by four attributes, UserID, MovieID, Rating and Timestamp. The attribute  
350 Timestamp has a different value for every record. We delete the attribute because it has almost no effects on clustering results.

The preprocessing process is described as follows: select objects that are in the range of  $x < -0.2, y > 0.4$  or  $x > 0.2, y > 0.4$  or  $x < -0.2, y < -0.5$  or  $x > 0.2, y < -0.3$  in the coordinate system as a new set after the visualization  
355 of the initial data set; then visualize the new data in Fig.4.

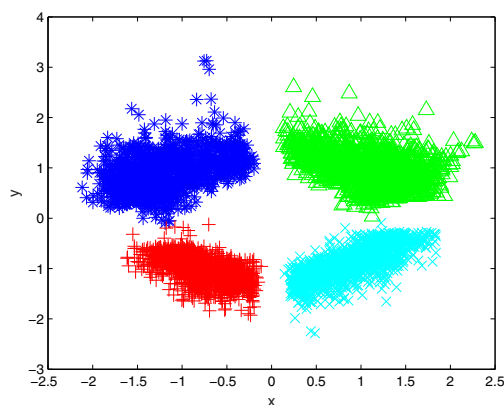


Figure 4: The distribution of MovieLens data after preprocessing.

<sup>3</sup><http://grouplens.org/datasets/movielens/>

Apparently, the data seen from Fig.4 are divided into 4 clusters.

The final data sets after preprocessing are listed in Table 4. These data sets are used to evaluate the  $k$ -mw-modes algorithm.

Table 4: Data sets after preprocessing.

Data set	Matrix-objects	Attributes	Records	$k$
Web	2401	2	22250	2
Basket	900	2	6300	3
Alibaba	793	2	165655	3
Musk	92	167	476	2
Movie	4592	3	728487	4

#### 4.2. Evaluation indexes

360 To evaluate the effectiveness of the  $k$ -mw-modes clustering algorithm, we used the following five external criteria: (1) adjusted rand index (ARI) [7], (2) normalized mutual information (NMI) [8], (3) accuracy (AC), (4) precision (PE) and (5) recall (RE) to measure the similarity between two partitions of objects in a given data set.

365 Let  $\mathbf{X}$  be a matrix-object data set,  $C = \{C_1, C_2, \dots, C_k\}$  be a clustering result of  $\mathbf{X}$ ,  $P = \{P_1, P_2, \dots, P_k\}$  be a real partition in  $\mathbf{X}$ . The overlap between  $C$  and  $P$  can be summarized in a contingency table shown in Table 5, where  $n_{ij}$  denotes the number of objects in common between  $P_i$  and  $C_j$ ,  $n_{ij} = |P_i \cap C_j|$ .  $p_i$  and  $c_j$  are the number of objects in  $P_i$  and  $C_j$ , respectively.

370 The five evaluation indexes are defined as follows:

$$ARI = \frac{\sum_{ij} C_{n_{ij}}^2 - [\sum_i C_{p_i}^2 \sum_j C_{c_j}^2] / C_n^2}{\frac{1}{2}[\sum_i C_{p_i}^2 + \sum_j C_{c_j}^2] - [\sum_i C_{p_i}^2 \sum_j C_{c_j}^2] / C_n^2},$$

$$NMI = \frac{\sum_{i=1}^k \sum_{j=1}^{k'} n_{ij} \log(\frac{n_{ij}n}{p_i c_j})}{\sqrt{\sum_{i=1}^k p_i \log(\frac{p_i}{n}) \sum_{j=1}^{k'} c_j \log(\frac{c_j}{n})}},$$

Table 5: The contingency table.

	$C_1$	$C_2$	.....	$\cdots$	$C_{k'}$	$Sums$
$P_1$	$n_{11}$	$n_{12}$		$\cdots$	$n_{1k'}$	$p_1$
$P_2$	$n_{21}$	$n_{22}$		$\cdots$	$n_{2k'}$	$p_2$
$\vdots$	$\vdots$	$\vdots$		$\ddots$	$\vdots$	$\vdots$
$P_k$	$n_{k1}$	$n_{k2}$		$\cdots$	$n_{kk'}$	$p_k$
$Sums$	$c_1$	$c_2$		$\cdots$	$c_{k'}$	$n$

$$AC = \frac{1}{n} \max_{j_1 j_2 \cdots j_k \in S} \sum_{i=1}^k n_{ij_i},$$

$$PE = \frac{1}{k} \sum_{i=1}^k \frac{n_{ij_i^*}}{p_i},$$

$$RE = \frac{1}{k'} \sum_{i=1}^{k'} \frac{n_{ij_i^*}}{c_i},$$

where  $n_{1j_1^*} + n_{2j_2^*} + \cdots + n_{kj_k^*} = \max_{j_1 j_2 \cdots j_k \in S} \sum_{i=1}^k n_{ij_i}$  ( $j_1^* j_2^* \cdots j_k^* \in S$ ) and  $S = \{j_1 j_2 \cdots j_k : j_1, j_2, \cdots, j_k \in \{1, 2, \cdots, k\}, j_i \neq j_t \text{ for } i \neq t\}$  is a set of all permutations of  $1, 2, \cdots, k$ . For  $AC, PE, RE$ ,  $k$  is equal to  $k'$  in general case. In addition, we consider that the higher the values of  $ARI, NMI, AC, PE$  and  $RE$  are, the better the clustering solution is.

#### 4.3. Comparisons of two cluster center update algorithms GAFMWM and HAFMWM

In this section, we show comparison results of two cluster center update algorithms *GAFMWM* and *HAFMWM* used in the  $k$ -mw-modes algorithm. As the  $k$ -mw-modes algorithm that is the extension of the  $k$ -modes algorithm is sensitive to the initial cluster centers, we executed the  $k$ -mw-modes algorithm with two cluster center update algorithms 10 times, respectively. And all of our experiments were conducted on a PC with an Intel Xeon CPU I7(3.4GHz) and 16GB memory. Table 6 shows the average values and standard deviations of clustering evaluation indexes on Market Basket data set.



Table 6: Comparison results of the  $k$ -mw-modes algorithm with  $GAFMWM$  and  $HAFMWM$  on Market Basket data.

	$AC$	$PE$	$RE$	$ARI$	$NMI$
$k$ -mw-modes+ $GAFMWM$	$0.9058 \pm 0.1176$	$0.9072 \pm 0.1151$	$0.9004 \pm 0.1262$	$0.8015 \pm 0.2215$	$0.8099 \pm 0.2022$
$k$ -mw-modes+ $HAFMWM$	$0.8834 \pm 0.1370$	$0.8906 \pm 0.1385$	$0.8726 \pm 0.1483$	$0.7383 \pm 0.2419$	$0.7180 \pm 0.2162$

385 In addition, we compared the execution time of the  $k$ -mw-modes algorithm with two update methods and the experimental results are shown in Table 7.

Table 7: Run-time of the  $k$ -mw-modes algorithm with  $GAFMWM$  and  $HAFMWM$  on Market Basket data.

	Run-time (Second)
$k$ -mw-modes + $GAFMWM$	$1.60497 \times 10^5 \pm 0.4731 \times 10^5$
$k$ -mw-modes + $HAFMWM$	$15.6146 \pm 4.2214$

From Table 6, we can find that the  $k$ -mw-modes algorithm with  $GAFMWM$  is better than it with  $HAFMWM$  on all evaluation indexes. But we also can see that the  $k$ -mw-modes algorithm with  $GAFMWM$  is very time consuming from  
 390 Table 7. It took about 45 hours to produce one clustering result on Market Basket data, which is not acceptable in real applications. However, the  $k$ -mw-modes algorithm with  $HAFMWM$  only took a few seconds to produce a clustering result on the same data set.  $HAFMWM$  compared with  $GAFMWM$  speeds up the  $k$ -mw-modes process tremendously. Furthermore, the values of  
 395 evaluation indexes of both algorithms are closer. Therefore, we will use the  $k$ -mw-modes algorithm with  $HAFMWM$  instead of  $GAFMWM$  in the following experiments.

#### 4.4. Comparisons of the $k$ -mw-modes with other algorithms

As far as we know, no appropriate algorithms can be used to directly cluster  
 400 the categorical matrix-object data. To cluster the data by some existed algorithms, we have to convert the representation of matrix-objects.

If we use a feature vector to represent a matrix-object, the matrix-object data can be clustered by  $k$ -modes-type algorithms. We choose the attribute value with maximal frequency as a representative on each attribute for each matrix-object. Then, each matrix-object is simplified a feature vector. In this subsection, we compared the  $k$ -mw-modes with the  $k$ -modes, the  $Wk$ -modes [16], the  $k$ -modes with the new dissimilarity measure (abbr.  $k$ -modes-cao) [10].

If we regard all categories to be independent, the matrix-object data can be identified with co-occurrence information data among objects and categories. Thus, some co-clustering algorithms [20, 21, 22] can be used to cluster the matrix-object data. For a given matrix-object data set, we count the number of each attribute value in each object to transform each object into a  $p$ -dimensional feature vector ( $p = \sum_{s=1}^m |V^s|$ ). In this subsection, we compared the  $k$ -mw-modes with FCCM that is one of co-clustering methods [20].

In this experiment, we clustered every data set 50 times respectively by the these algorithms and  $\varepsilon$  was set to 0.1 in the  $k$ -mw-modes algorithm. In the FCCM, we set  $T_u = 0.1$ ,  $T_w = 1.5$  and  $\epsilon = 0.0001$  [20]. The results of the five data sets are listed in Tables 8-12, respectively.

Table 8: Comparison results of the five algorithms on Microsoft Web data.

<i>Algorithms</i>	<i>AC</i>	<i>PE</i>	<i>RE</i>	<i>ARI</i>	<i>NMI</i>
$k$ -modes	0.6117±0.0608	0.6101±0.0598	0.5827±0.0802	0.0592±0.0613	0.0414±0.0432
$Wk$ -modes	0.6284±0.0742	0.7158±0.0615	0.6009±0.1002	0.0806±0.0876	0.0852±0.0786
$k$ -modes-cao	0.6504±0.0653	0.6521±0.0690	0.6357±0.0789	0.1048±0.0835	0.0790±0.0654
FCCM					
$k$ -MW-modes	0.9220±0.0940	0.9288±0.0962	0.9148±0.1013	0.7463±0.1918	0.6806±0.1745

The left of the symbol “±” in these tables represents the means of external criterion values for 50 experiments and the right of it represents standard deviations. From the five tables, we can find that the means on AC, PE, RE, ARI, NMI in the  $k$ -mw-modes algorithm are generally higher than that in the  $k$ -modes, the  $Wk$ -modes and the  $k$ -modes-cao algorithms. What’s more, the  $k$ -mw-modes algorithm on the index AC is approximately 10%~30% more than

Table 9: Comparison results of the five algorithms on Market Basket data.

<i>Algorithms</i>	<i>AC</i>	<i>PE</i>	<i>RE</i>	<i>ARI</i>	<i>NMI</i>
<i>k</i> -modes	0.4939±0.0670	0.4967±0.0723	0.4865±0.0883	0.0836±0.0564	0.0821±0.0535
<i>Wk</i> -modes	0.5213±0.0643	0.6375±0.0815	0.5240±0.1004	0.1095±0.0615	0.1620±0.0638
<i>k</i> -modes-cao	0.5478±0.0605	0.5806±0.0763	0.5449±0.0738	0.1373±0.0520	0.1517±0.0529
FCCM	0.5405±0.0211	0.5629±0.0458	0.6711±0.1167	0.1522±0.0205	0.1690±0.0233
<i>k</i> -MW-modes	0.8834±0.1370	0.8906±0.1385	0.8726±0.1483	0.7383±0.2419	0.7180±0.2162

Table 10: Comparison results of the five algorithms on Alibaba data.

<i>Algorithms</i>	<i>AC</i>	<i>PE</i>	<i>RE</i>	<i>ARI</i>	<i>NMI</i>
<i>k</i> -modes	0.4817±0.0000	0.4817±0.0012	0.3333±0.0000	0.0000±0.0020	0.0022±0.0015
<i>Wk</i> -modes	0.4832±0.0011	0.7796±0.0821	0.3353±0.0015	-0.0012±0.0040	0.0217±0.0085
<i>k</i> -modes-cao	0.4817±0.0000	0.4805±0.0024	0.3333±0.0000	0.0036±0.0040	0.0049±0.0030
FCCM	0.6528±0.0459	0.6217±0.0564	0.5840±0.0525	0.2866±0.0732	0.2245±0.0610
<i>k</i> -MW-modes	0.6528±0.0459	0.6217±0.0564	0.5840±0.0525	0.2866±0.0732	0.2245±0.0610

Table 11: Comparison results of the five algorithms on Musk data.

<i>Algorithms</i>	<i>AC</i>	<i>PE</i>	<i>RE</i>	<i>ARI</i>	<i>NMI</i>
<i>k</i> -modes	0.5567±0.0382	0.5639±0.0434	0.5547±0.0398	0.0089±0.0239	0.0155±0.0192
<i>Wk</i> -modes	0.5159±0.0070	0.6141±0.0988	0.5070±0.0107	-0.0010±0.0014	0.0182±0.0231
<i>k</i> -modes-cao	0.5387±0.0259	0.5402±0.0278	0.5361±0.0295	-0.0021±0.0111	0.0067±0.0086
FCCM	0.5846±0.0285	0.5976±0.0398	0.5820±0.0272	0.0224±0.0206	0.0285±0.0200
<i>k</i> -MW-modes	0.5600±0.0413	0.5688±0.0487	0.5565±0.0423	0.0116±0.0276	0.0183±0.0243

Table 12: Comparison results of the five algorithms on MovieLens data.

<i>Algorithms</i>	<i>AC</i>	<i>PE</i>	<i>RE</i>	<i>ARI</i>	<i>NMI</i>
<i>k</i> -modes	0.4885±0.1206	0.5191±0.1508	0.4864±0.1194	0.1997±0.1205	0.2360±0.1416
<i>Wk</i> -modes	0.5267±0.1220	0.6147±0.1235	0.5351±0.1248	0.2617±0.1231	0.3251±0.1474
<i>k</i> -modes-cao	0.4348±0.0809	0.4266±0.0733	0.4311±0.0794	0.1301±0.0809	0.1261±0.0754
FCCM	0.6244±0.0303	0.6851±0.0297	0.6233±0.0321	0.3211±0.0424	0.3904±0.0367
<i>k</i> -MW-modes	0.6244±0.0303	0.6851±0.0297	0.6233±0.0321	0.3211±0.0424	0.3904±0.0367

425 those algorithms on it except for Musk data. We can also see that the AC in the *Wk*-modes and *k*-modes-cao algorithms is generally higher than it in the

$k$ -modes algorithm.

For the FCCM, its accuracy on Musk data is higher than that of the proposed algorithm, but it has not clustering results on Alibaba data and MovieLens data. The reason is that we cannot compute the memberships of attribute values for each cluster after normalization because some attribute values have a fairly high frequency in each object and the membership of every attribute value in a cluster defined in the FCCM increases exponentially as the increasing of the total frequency of the value in each object. For Microsoft Web data, we only can obtain a cluster, which is meaningless in real applications. Therefore, the FCCM is not suitable for the clustering of this type of data.

In short, it can be seen that the  $k$ -mw-modes algorithm is exactly better than the other four algorithms.

#### 4.5. Impact of $\varepsilon$

In the  $k$ -mw-modes algorithm, the parameter  $\varepsilon$  is used to determine whether the algorithm stops or not. How to decide the size of  $\varepsilon$  is a very difficult problem, because the clustering results may have some differences when different values of  $\varepsilon$  are selected. To analyze the impact of  $\varepsilon$  on the clustering performance of the  $k$ -mw-modes algorithm, we ran the algorithm 30 times with different values of  $\varepsilon$ , from 0.02 to 0.2 with step 0.01, on the five data sets and recorded the means of accuracy (AC) and iterations. Results of AC and iterations on the five data sets are shown in Fig.5 and Fig.6, respectively.

From Fig. 5, we can observe that the AC on Musk, Alibaba and MovienLens has almost no change and the AC on the other two data has been fluctuating as the increasing of  $\varepsilon$ , but overall they are relatively stable. Meanwhile, from Fig.6, we can see clearly that the iterations on the five data all show a decreasing trend on the whole as the increasing of  $\varepsilon$ . Particularly, for Web and Alibaba, the iterations decrease rapidly when  $\varepsilon < 0.1$  and they decrease slowly relatively when  $\varepsilon > 0.1$ . Higher AC and fewer iterations are expected in this algorithm, so we set  $\varepsilon = 0.1$ .

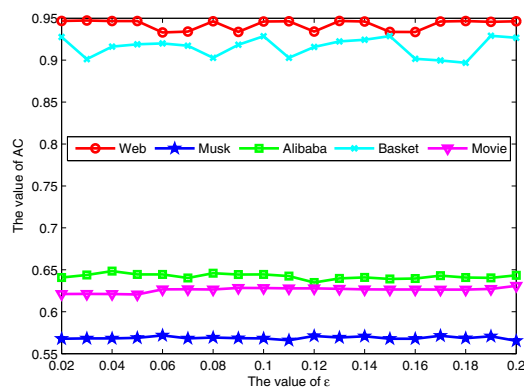


Figure 5: The AC of the five data with different  $\epsilon$ .

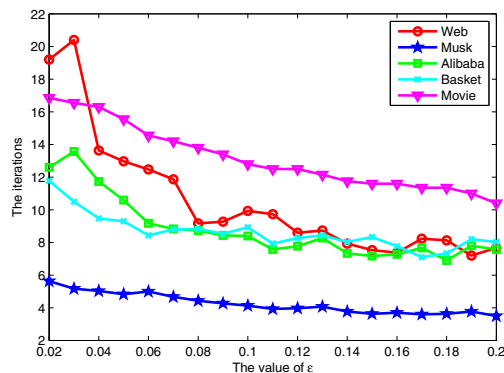


Figure 6: The iterations of the five data with different  $\epsilon$ .

## 5. Related work

In real applications, categorical data are widespread. The  $k$ -modes algorithm [3] extends the  $k$ -means algorithm [4] by using a simple matching dissimilarity measure for categorical objects, modes instead of means for clusters, and a frequency-based method to update modes in the clustering process to minimize the clustering objective function. These extensions have removed the numeric-only limitation of the  $k$ -means algorithm and enable the  $k$ -means clustering process to be used to efficiently cluster large categorical data sets from real

world [9, 10]. The  $k$ -modes clustering algorithm [11, 12] was also proposed  
465 independently. Huang [13] gave the relationship of the two  $k$ -modes methods.  
So far, the  $k$ -modes algorithm and its variants [14, 15, 16], including the fuzzy  
 $k$ -modes algorithm [17], the fuzzy  $k$ -modes algorithm with fuzzy centroid [18],  
the  $k$ -prototype algorithm [3] and the  $w$ - $k$ -means [19] have been used widely in  
470 many domains. However, these methods can not cluster matrix-object data set  
effectively.

## 6. Conclusions

In many database applications, the behavioural traits of a customer are car-  
ried in a detail table instead of a master table. To find the customer groups  
with different behavioural traits, a  $k$ -mw-modes algorithm was proposed for clus-  
475 tering categorical matrix-object data. In the proposed algorithm, the distance  
between two matrix-objects was defined and the representation and update ways  
of cluster centers were developed further. The convergence of the proposed al-  
gorithm was proved and the corresponding time complexity was analyzed as  
well. To speed up the clustering process, a heuristic method was proposed to  
480 construct multi-weighted-modes centers in each iteration of the  $k$ -mw-modes  
algorithm. Experimental results on the five real data sets have shown that  
the  $k$ -mw-modes algorithm is better than the  $k$ -modes-type algorithms and the  
FCCM in clustering categorical matrix-object data.

## Acknowledgements

485 This work was supported by the National Natural Science Foundation of  
China (under grants 61573229, 61473194, 61432011 and U1435212), the Nat-  
ural Science Foundation of Shanxi Province (under grant 2015011048), the  
Shanxi Scholarship Council of China (under grant 2016-003) and the National  
Key Basic Research and Development Program of China (973) (under grant  
490 2013CB329404).

**References**

- [1] J. Han, M. Kamber, J. Pei, Data Mining: Concepts and Techniques, 3rd Edition, Morgan Kaufmann Publishers Inc., San Francisco, CA, USA, 2011.
- [2] R. Xu, D. Wunsch, Clustering, Vol. 10, John Wiley & Sons, 2008.
- 495 [3] Z. Huang, Extensions to the  $k$ -means algorithm for clustering large data sets with categorical values, Data Mining and Knowledge Discovery 2 (3) (1998) 283–304.
- [4] J. MacQueen, Some methods for classification and analysis of multivariate observations, in: Proceedings of the fifth Berkeley Symposium on Mathematical Statistics and Probability, Vol. 1, California, USA, 1967, pp. 281–  
500 297.
- [5] S. Schiffman, L. Reynolds, F. Young, Introduction to Multidimensional Scaling: Theory, Methods, and Applications, Academic Press, 1981.
- [6] K. Bache, M. Lichman, UCI machine learning repository (2014).  
505 URL <http://archive.ics.uci.edu/ml>
- [7] J. Liang, L. Bai, C. Dang, F. Cao, The  $k$ -means type algorithms versus imbalanced data distributions, IEEE Transactions on Fuzzy Systems 20 (4) (2012) 728–745.
- [8] A. Strehl, J. Ghosh, Cluster ensembles—a knowledge reuse framework for  
510 combining multiple partitions, The Journal of Machine Learning Research 3 (2003) 583–617.
- [9] F. Cao, J. Z. Huang, J. Liang, Trend analysis of categorical data streams with a concept change method, Information Sciences 276 (2014) 160–173.
- [10] F. Cao, J. Liang, D. Li, L. Bai, C. Dang, A dissimilarity measure for the  $k$ -  
515 modes clustering algorithm, Knowledge-Based Systems 26 (2012) 120–127.

- [11] J. D. Carroll, A. Chaturvedi, P. E. Green, *k*-means, *k*-medians and *k*-modes: Special cases of partitioning multiway data., 1994.
- [12] A. Chaturvedi, P. E. Green, J. D. Carroll, *k*-modes clustering, *Journal of Classification* 18 (1) (2001) 35–55.
- 520 [13] Z. Huang, M. K. Ng, A note on *k*-modes clustering, *Journal of Classification* 20 (2) (2003) 257–261.
- [14] M. K. Ng, M. J. Li, J. Z. Huang, Z. He, On the impact of dissimilarity measure in *k*-modes clustering algorithm, *IEEE Transactions on Pattern Analysis and Machine Intelligence* 29 (3) (2007) 503–507.
- 525 [15] L. Bai, J. Liang, C. Dang, F. Cao, The impact of cluster representatives on the convergence of the *k*-modes type clustering, *IEEE Transactions on Pattern Analysis and Machine Intelligence* 35 (6) (2013) 1509–1522.
- [16] F. Cao, J. Liang, D. Li, X. Zhao, A weighting *k*-modes algorithm for subspace clustering of categorical data, *Neurocomputing* 108 (2013) 23–30.
- 530 [17] Z. Huang, M. K. Ng, A fuzzy *k*-modes algorithm for clustering categorical data, *IEEE Transactions on Fuzzy Systems* 7 (4) (1999) 446–452.
- [18] D.-W. Kim, K. H. Lee, D. Lee, Fuzzy clustering of categorical data using fuzzy centroids, *Pattern Recognition Letters* 25 (11) (2004) 1263–1271.
- [19] J. Z. Huang, M. K. Ng, H. Rong, Z. Li, Automated variable weighting  
535 in *k*-means type clustering, *IEEE Transactions on Pattern Analysis and Machine Intelligence* 27 (5) (2005) 657–668.
- [20] C.-H. Oh, K. Honda, H. Ichihashi, Fuzzy clustering for categorical multivariate data, *Proc. of Joint 9th IFSA World Congress and 20th NAFIPS Int. Conf.*, (2001) 2154–2159.
- 540 [21] K. Kummamuru, A. Dhawale, R. Krishnapuram, Fuzzy co-clustering of documents and keywords, *Proc. 2003 IEEE Int. Conf. Fuzzy Systems*, 2 (2003) 772–777.



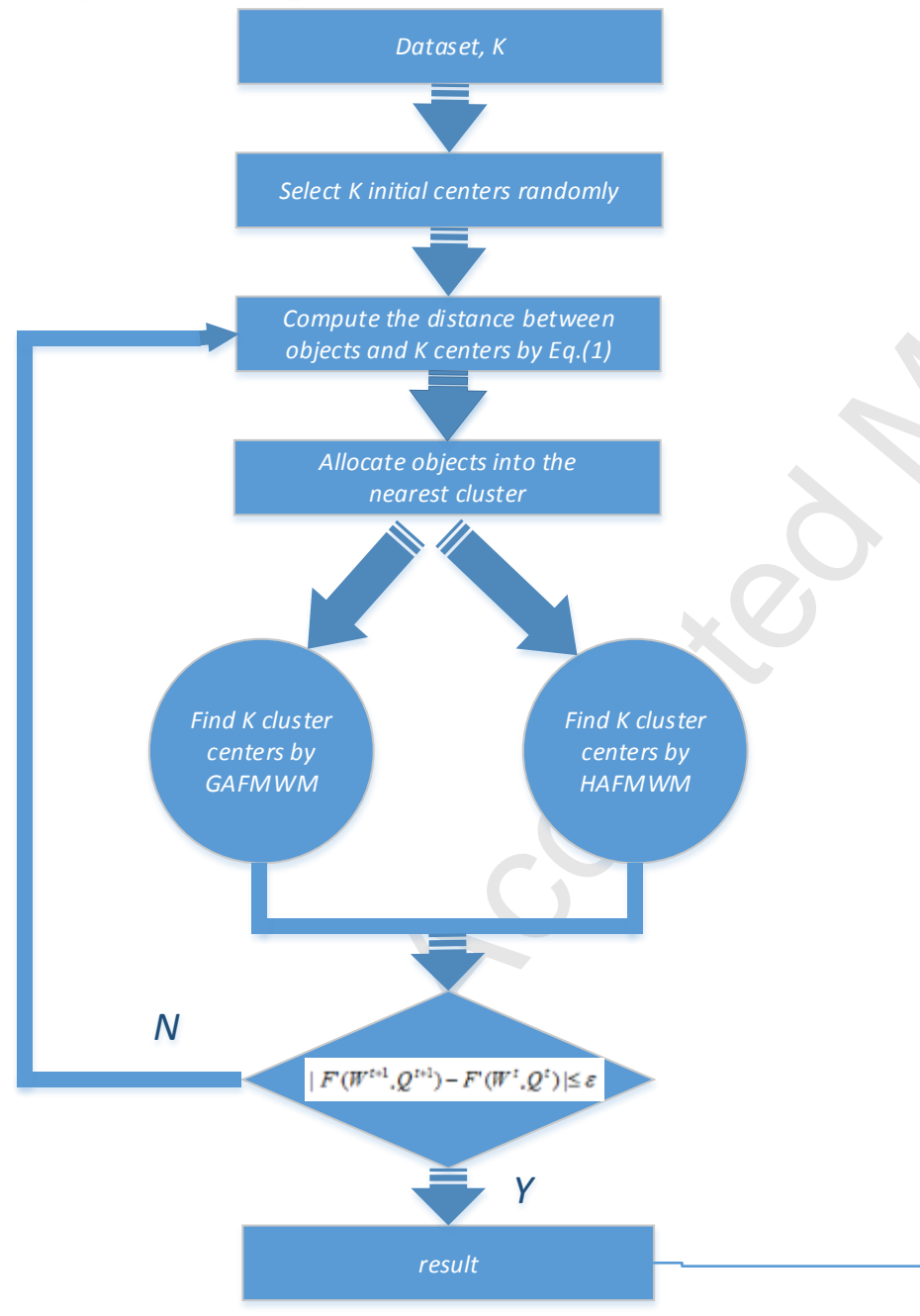
- [22] H. Frigui, O. Nasraoui, Simultaneous categorization of text documents and identification of cluster dependent keywords, Proc. 2002 IEEE Int. Conf. Fuzzy Systems, 2 (2002) 1108-1113.

545

Accepted Manuscript

- (1) We propose a k-multi-weighted-modes (abbr. k-mw-modes) algorithm for clustering categorical matrix-object data and the k-modes algorithm is its special case.
- (2) We give a heuristic method to choose the locally optimal Multi-Weighted-modes in the iteration of the k-mw-modes algorithm and the update process of the k-modes algorithm is its special case.
- (3) Experimental results on the five real data sets from different applications have shown the effectiveness of the k-mw-modes algorithm.

The proposed algorithm (k-mw-modes)



Categorical matrix-object data

Table 1: A real application example from Taobao.

User_ID	Brand_ID	Brand_Name	Visit_Date	Visited_Times
10944750	13451	WETHERM	06-04	8
	21110	SEMIR	06-07	1
	25687	JOSINY	05-08	11
	25687	JOSINY	05-15	4
	25687	JOSINY	05-17	2
	25687	JOSINY	06-06	2
	25687	JOSINY	06-15	2
	25687	JOSINY	07-02	5
	25687	JOSINY	07-25	2
	25687	JOSINY	08-09	3
25687	JOSINY	08-13	1	
8149250	18805	YEAR CON	05-25	1
	18805	YEAR CON	05-26	4
	18805	YEAR CON	06-12	1
	13451	SEMIR	06-30	1
...	...	...	...	...
1832000	25687	JOSINY	08-14	18
	25687	JOSINY	08-15	3
	25687	JOSINY	06-26	5
	25687	JOSINY	08-11	1
	21110	OBERORA	08-03	1
	10151	YEAR CON	07-26	5

