# Categorical data visualization and clustering using subjective factors

## Chia-Hui Chang *, Zhi-Kai Ding

*Department of Computer Science and Information Engineering, National Central University, No. 300,
Jhungda Road, Jhungli City, Taoyuan 320, Taiwan*

## Abstract

Clustering is an important data mining problem. However, most earlier work on clustering focused on numeric attributes which have a natural ordering to their attribute values. Recently, clustering data with categorical attributes, whose attribute values do not have a natural ordering, has received more attention. A common issue in cluster analysis is that there is no single correct answer to the number of clusters, since cluster analysis involves human subjective judgement. Interactive visualization is one of the methods where users can decide a proper clustering parameters. In this paper, a new clustering approach called CDCS (Categorical Data Clustering with Subjective factors) is introduced, where a visualization tool for clustered categorical data is developed such that the result of adjusting parameters is instantly reflected. The experiment shows that CDCS generates high quality clusters compared to other typical algorithms.
© 2004 Published by Elsevier B.V.

*Keywords:* Data mining; Cluster analysis; Categorical data; Cluster visualization

---

* Corresponding author. Fax: +886 3 4222681.
*E-mail addresses:* chia@csie.ncu.edu.tw (C.-H. Chang), sting@db.csie.ncu.edu.tw (Z.-K. Ding).

## 1. Introduction

Clustering is one of the most useful tasks for discovering groups and identifying interesting distributions and patterns in the underlying data. The clustering problem is about partitioning a given data set into groups (clusters) such that the data points in a cluster are more similar to each other than points in different clusters. The clusters thus discovered are then used for describing characteristics of the data set. Cluster analysis has been widely used in numerous applications, including pattern recognition, image processing, land planning [21], text query interface [11], market research, etc.

Many clustering methods have been proposed in the literature and most of these handle data sets with numeric attributes, where proximity measure can be defined by geometrical distance. For categorical data which has no order relationships, a general method is to transform it into binary data. However, such binary mapping may lose the meaning of original data set and result in incorrect clustering, as reported in [7]. Furthermore, high dimensions will require more space and time if the similarity function involves with matrix computation, such as the Mahalanobis measure [16].

Another problem we face in clustering is how to validate the clustering results and decide the optimal number of clusters that fits a data set. Most clustering algorithms require some predefined parameters for partitioning. These parameters influence the clustering result directly. For a specific application, it may be important to have well separated clusters, while for another it may be more important to consider the compactness of the clusters. Hence, there is no correct answer for the optimal number of clusters since cluster analysis may involve human subjective judgement, and visualization is one of the most intuitive ways for users to decide a proper clustering.

In Fig. 1 for example, there are 54 objects displayed. For some people, there are two clusters, while some may think there are six clusters, still others may think there are 18 clusters, depending on their subjective judgement. In other words, a value can be small in a macroscopic view, but it can be large in a microscopic view. The definition of similarity varies with respect to different views. Therefore, if categorical data can be visualized, parameter adjustment can be easily done, even if several parameters are involved.

In this paper, we present a method for visualization of the clustered categorical data such that users' subjective factors can be reflected by adjusting clustering parameters, and therefore to increase the clustering result's reliability. The proposed method, CDCS (Categorical Data Clustering using Subjective factors), can be divided into three steps: (1) the first step incorporates a single-pass clustering method to group objects with high similarity, (2) then, small clusters are merged
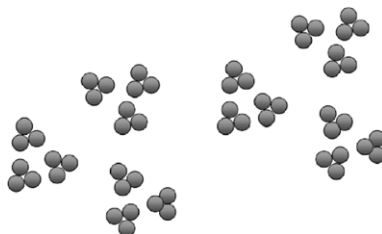
Fig. 1. Fifty-four objects displayed in a plane.

55  and displayed for visualization, and (3) through the proposed interactive visualization tool, users
56  can observe the data set and determine appropriate parameters for clustering.
57      The rest of the paper is organized as follows. Section 2 reviews related work in categorical data
58  clustering and data visualization. Section 3 introduces the architecture of CDCS and the cluster-
59  ing algorithm utilized. Section 4 discusses the visualization method of CDCS in detail. Section 5
60  presents an experimental evaluation of CDCS using popular data sets and comparisons with two
61  famous algorithms AutoClass [2] and k-mode [9]. Section 6 presents the conclusions and suggests
62  future work.

## 63  2. Related work

64      Clustering is broadly recognized as a useful tool for many applications. Researchers of many
65  disciplines (such as databases, machine learning, pattern recognition, statistics, etc.) have ad-
66  dressed clustering problem in many ways. However, most researches concern numerical type data
67  which has geometrical shape and clear distance definition, while little attention has been paid to
68  categorical data clustering. Visualization is an interactive, reflective method that supports explo-
69  ration of data sets by dynamically adjusting parameters to see how they affect the information
70  being presented. In this section, we review works on categorical data clustering and visualization
71  methods.

### 72  2.1. Categorical data clustering

73      In recent years, a number of clustering algorithms for categorical data have been proposed,
74  partly due to increasing applications in market baskets analysis, customer databases, etc. We
75  briefly review the main algorithms below.
76      One of the most common ways to solve categorical data clustering is to extend existing algo-
77  rithms with a proximity measure for categorical data, and many clustering algorithms belong
78  to this category. For example, k-mode [9] is based on k-mean [15] but adopts a new similarity
79  function to handle categorical data. A cluster center for k-mode is represented by a virtual object
80  with the most frequent attribute values in the cluster. Thus, the distance between two data objects
81  is defined as the number of different attribute values.
82      ROCK [7] is a bottom-up clustering algorithm which adopts a similarity function based on the
83  number of common neighbors, which is defined by the Jaccard coefficient [18]. Two data objects
84  are more similar if they have more common neighbors. Since the time complexity for the bottom-
85  up hierarchical algorithm is quadratic, it clusters a randomly sampled data set and then partitions
86  the entire data set based on these clusters. COBWEB [3], on the other hand, is a top-down clus-
87  tering algorithm which constructs a classification tree to record cluster information. The disad-
88  vantage of COBWEB is that its classification tree is not height-balanced for skewed input data,
89  which may cause increased time and space cost.
90      AutoClass [2] is an EM-based approach which is supplemented by a Bayesian evaluation for
91  determining the optimal classes. It assumes a predefined distribution for data and tries to maxi-
92  mize the function with appropriate parameters. AutoClass requires a range for the number of

93  clusters as input. Because AutoClass is a typical iterative clustering algorithm, if the data cannot
94  be entirely loaded into memory, the time cost will be expensive.
95      STIRR [6] is an iterative method based on non-linear dynamical systems. Instead of clustering
96  objects themselves, the algorithm aims at clustering co-occur attribute values. Usually, this ap-
97  proach can discover the largest cluster with the most attribute values, and even if the idea of
98  orthonormal is introduced [23], it can only discover one more cluster. CACTUS [5] is another
99  algorithm that conducts clustering from attribute relationship. CACTUS employs a combination
100 of inter-attribute and intra-attribute summaries to find clusters. However, there has been no re-
101 port on how such an approach can be used for clustering general data sets.
102     Finally, several endeavors have been tried to mine clusters with association rules. For example,
103 Kosters et al. proposed the clustering of a specific type of data sets where the objects are vectors of
104 binary attributes from association rules [8]. The hypergraph-based clustering in [13] is used to par-
105 tition items and then transactions based on frequent itemsets. Wang et al. [22] also focus on trans-
106 action clustering as in [8].

107 *2.2. Visualization methods*

108     Since human vision is endowed with the classification ability for graphic figures, it would
109 greatly help solving the problem if the data could be graphically transformed for visualization.
110 However, human vision is only useful for figures with low dimension. For high-dimensional data,
111 it must be mapped into low dimensions for visualization, and there are several visualization meth-
112 ods including linear mapping and non-linear mapping. Linear mapping, like principle component
113 analysis, is effective but cannot truly reflect the data structure. Non-linear mapping, like Sammon
114 projection [12] and SOM [10] requires more computation but is better at preserving the data struc-
115 ture. However, whichever method is used, traditional visualization methods can transform only
116 numerical data. For categorical data, visualization is only useful for attribute dependency analysis
117 and is not helpful for clustering. For example, the mosaic display method [4], a popular statistical
118 visualization method, displays the relationships between two attributes. Users can view the dis-
119 play in a rectangle composed of many mosaic graphs and compare it to another mosaic display,
120 assuming that two attributes are independent. The tree map [19], which is the only visualization
121 method which can be used for distribution analysis, transforms data distribution to a tree com-
122 posed of many nodes. Each node in this tree is displayed as a rectangle and its size represents
123 the frequency of an attribute value. Users can thus observe an overview of the attribute distribu-
124 tion. However, it provides no insight to clustering. Finally, reordering of attribute values may help
125 visualization for categorical data with one attribute as proposed in [14].

126 **3. Categorical data clustering and visualization**

127     In this section, we introduce the clustering and visualization approach in our framework. We
128 utilize the concept of Bayesian classifiers as a proximity measure for categorical data. The process
129 of CDCS can be divided into three steps: (1) in the first step, it applies "simple clustering seeking"
130 [20] to group objects with high similarity, (2) then, small clusters are merged and displayed by cat-
131 egorical cluster visualization, (3) users can then adjust the merging parameters and view the result
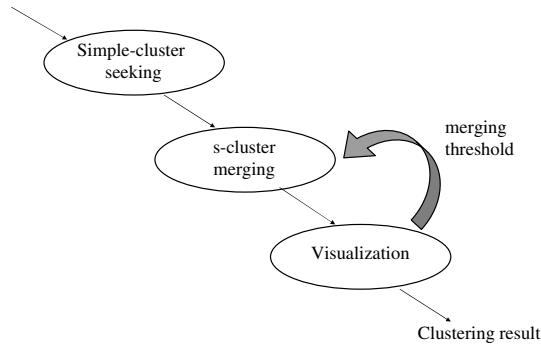
Fig. 2. The CDCS process.

132 through the interactive visualization tool. The process continues until users are satisfied with the
133 result (see Fig. 2).

134    Simple cluster seeking, sometimes called dynamic clustering, is a one pass clustering algorithm
135 which does not require the specification of the desired number of clusters. Instead, a similarity
136 threshold is used to decide if a data object should be grouped into an existing cluster or form a
137 new cluster. More specifically, the data objects are processed individually and sequentially. The
138 first data object forms a single cluster by itself. Next, each data object is compared to existing clus-
139 ters. If its similarity with the most similar cluster is greater than a given threshold, this data object
140 is assigned to that cluster and the representation of that cluster is updated. Otherwise, a new clus-
141 ter is formed. The advantage of dynamic clustering is that it provides simple and incremental clus-
142 tering where each data sample contributes to changes in the clusters. Besides, the time complexity,
143 $O(kn)$, for clustering $n$ objects into $k$ clusters is suitable for handling large data sets.

144    However, there is one inherent problem for this dynamic clustering: the clustering result can be
145 affected by the input order of data objects. To solve this problem, higher similarity thresholds can
146 be used to decrease the influence of the data order and ensure that only highly similar data objects
147 are grouped together. As a large number of small clusters (called s-clusters) can be produced, the
148 cluster merging step is required to group s-clusters into larger groups. Therefore, a merging step
149 similarity threshold is designed to be adjusted for interactive visualization. Thus, a user's views
150 about the clustering result can be extracted when he/she decides a proper threshold.

151 *3.1. Proximity measure for categorical data*

152    Clustering is commonly known as an unsupervised learning process. The simple cluster seeking
153 approach can be viewed as a classification problem since it predicts whether a data object belongs
154 to an existing cluster or class. In other words, data in the same cluster can be considered as having
155 the same class label. Therefore, the similarity function of a data object to a cluster can be repre-
156 sented by the probability that the data object belongs to that cluster. Here, we adopt a similarity
157 function based on the naive Bayesian classifier [17], which is used to compute the largest posteriori
158 probability $\text{Max}_j P(C_j|X)$ for a data object $X = (v_1, v_2, \ldots, v_d)$ to an existing cluster $C_j$. Using
159 Bayes' theorem, $P(C_j|X)$ can be computed by

*C.-H. Chang, Z.-K. Ding / Data & Knowledge Engineering xxx (2004) xxx–xxx*

$$P(C_j|X) \propto P(X|C_j)P(C_j) \tag{1}$$

162 Assuming attributes are conditionally independent, we can replace $P(X|C_j)$ by $\prod_{i=1}^{d}P(v_i|C_j)$,
163 where $v_i$ is $X$'s attribute value for the $i$th attribute. $P(v_i|C_j)$, a simpler form for $P(A_i = v_i|C_j)$, is
164 the probability of $v_i$ for the $i$th attribute in cluster $C_j$, and $P(C_j)$ is the priori probability defined
165 as the number of objects in $C_j$ to the total number of objects observed.
166 Applying this idea in dynamic clustering, the proximity measure of an incoming object $X_i$ to an
167 existing cluster $C_j$ can be computed as described above, where the prior objects $X_1, \ldots, X_{i-1}$ be-
168 fore $X_i$ are considered as the training set and objects in the same cluster are regarded as having the
169 same class label. For the cluster $C_k$ with the largest posteriori probability, if the similarity is great-
170 er than a threshold $g$ defined as

$$g = p^{d-e} \times \epsilon^e \times P(C_k) \tag{2}$$

173 then $X_i$ is assigned to cluster $C_k$ and $P(v_i|C_k)$, $i = 1, \ldots, d$, are updated accordingly. For each clus-
174 ter, a table is maintained to record the pairs of attribute value and their frequency for each attri-
175 bute. Therefore, to update $P(v_i|C_k)$ is simply an increase of the frequency count. Note that to
176 avoid zero product, the probability $P(v_i|C_j)$ is computed by $\frac{N_j(v_i)+m\cdot r}{|C_j|+m}$, where $N_j(v_i)$ is the number
177 of examples in cluster $C_j$ having attribute value $v_i$, $r$ is the reciprocal of the number of values
178 for the $i$th attribute as suggested in [17].
179 The equation for the similarity threshold is similar to the posteriori probability
180 $P(C_j|X) = \prod_{i=1}^{d}P(v_i|C_j)P(C_j)$, where the symbol $p$ denotes the average proportion of the highest
181 attribute value for each attribute, and $e$ denotes the number of attributes that can be allowed/tol-
182 erated for various values. For such attributes, the highest proportion of different attribute values
183 is given a small value $\epsilon$. This is based on the idea that the objects in the same cluster should possess
184 the same attribute values for most attributes, while some attributes may be quite dissimilar. In a
185 way, this step can be considered as density based clustering since probabilistic proximity measure
186 is the basis in the density based clustering. For large $p$ and small $e$, we will have many compact s-
187 clusters. In the most extreme situation, where $p = 1$ and $e = 0$, each distinct object is classified to a
188 cluster. CDCS adopts a default value 0.9 and 1 for $p$ and $e$, respectively. The resulting clusters are
189 usually small, highly condensed and applicable for most data sets.

190 *3.2. Group merging*

191 In the second step, we group the resulting s-clusters from dynamic clustering into larger clusters
192 ready for display with the proposed visualization tool. To merge s-clusters, we first compute the
193 similarity scores for each cluster pair. The similarity score between two s-clusters $C_x$ and $C_y$ is de-
194 fined as follows:

$$\text{sim}(C_x, C_y) = \prod_{i=1}^{d} \left[ \sum_{j}^{|A_i|} \min\{P(v_{i,j}|C_x), P(v_{i,j}|C_y)\} + \epsilon \right] \tag{3}$$

198 where $P(v_{i,j}|C_x)$ denotes the probability of the $j$th attribute value for the $i$th attribute in cluster $C_x$,
199 and $|A_i|$ denotes the number of attribute values for the $i$th attribute. The idea behind this defini-
200 tion is that the more the clusters intersect, the more similar they are. If the distribution of attribute

|   | 1 | 2 | 3 | 4 | 5 | 6 |
|---|---|---|---|---|---|---|
| 1 | 1 | 0 | 0 | 0 | 1 | 0 |
| 2 | 0 | 1 | 0 | 0 | 0 | 0 |
| 3 | 0 | 0 | 1 | 1 | 0 | 0 |
| 4 | 0 | 0 | 1 | 1 | 0 | 0 |
| 5 | 1 | 0 | 0 | 0 | 1 | 1 |
| 6 | 0 | 0 | 0 | 0 | 1 | 1 |

Resulting groups:
{1,5,6},
{2},
{3,4}

Fig. 3. Binary similarity matrix (BM).

values for two clusters is similar, they will have a higher similarity score. There is also a merge threshold $g'$, which is defined as follows:

$$g' = (p')^{d-e'} \times \epsilon^{e'} \tag{4}$$

Similar to the last section, the similarity threshold $g'$ is defined by $p'$, the average percentage of common attribute values for an attribute; and $e'$, the number of attributes that can be allowed/tolerated for various values. The small value $\epsilon$ is given to be the reciprocal of the number of samples in the data set.

For each cluster pair $C_x$ and $C_y$, the similarity score is computed and recorded in a $n \times n$ matrix SM, where $n$ is the number of s-clusters. Given the matrix SM and a similarity threshold $g'$, we compute a binary matrix BM (of size $n \times n$) as follows. If SM$[x, y]$ is greater than the similarity threshold $g'$, cluster $C_x$ and $C_y$ are considered similar and BM$[x, y] = 1$. Otherwise, they are dissimilar and BM$[x, y] = 0$. Note that the similarity matrix, SM, is computed only once after the single-pass clustering. For each parameter adjustment ($g'$) by the user, the binary matrix BM is computed without recomputing SM. Unless the parameter for the first step is changed, there is no need to recompute SM.

With the binary matrix BM, we then apply a transitive concept to group s-clusters. To illustrate this, in Fig. 3, clusters 1, 5, and 6 can be grouped in one cluster since clusters 1 and 5 are similar, and clusters 5 and 6 are also similar (the other two clusters are {2} and {3, 4}). This merging step requires O($n^2$) computation, which is similar to hierarchical clustering. However, the computation is conducted for $n$ s-clusters instead of data objects. In addition, this transitive concept allows arbitrarily shaped clusters to be discovered.

## 4. Visualization with CDCS

Simply speaking, visualization in CDCS is implemented by transforming a cluster into a graphic line connected by 3D points. These three dimensions represent the attributes, attribute values and the percentages of an attribute value in the cluster. These lines can then be observed in 3D space through rotations to see if they are close to each other. In the following, we first introduce the principle behind our visualization method; and then describe how it can help determine a proper clustering.

*C.-H. Chang, Z.-K. Ding / Data & Knowledge Engineering xxx (2004) xxx–xxx*

231 *4.1. Principle of visualization*

232     Ideally, each attribute $A_i$ of a cluster $C_x$ has an obvious attribute value $v_{i,k}$ such that the proba-
233 bility of the attribute value in the cluster, $P(A_i = v_{i,k}|C_x)$, is maximum and close to 100%. Therefore,
234 a cluster can be represented by these attribute values. Consider the following coordinate system
235 where the $X$ coordinate axis represents the attributes, the $Y$-axis represents attribute values corre-
236 sponding to respective attributes, and the $Z$-axis represents the probability that an attribute value is
237 in a cluster. Note that for different attributes, the $Y$-axis represents different attribute value sets. In
238 this coordinate system, we can denote a cluster by a list of $d$ 3D coordinates, $(i, v_{i,k}, P(v_{i,k}|C_x))$,
239 $i = 1, \ldots, d$, where $d$ denotes the number of attributes in the data set. Connecting these $d$ points,
240 we get a graphic line in 3D. Different clusters can then be displayed in 3D space to observe their
241 closeness.
242     This method, which presents only attribute values with the highest proportions, simplifies the
243 visualization of a cluster. Through operations like rotation or up/down movement, users can then
244 observe the closeness of s-clusters from various angles and decide whether or not they should be
245 grouped in one cluster. Graphic presentation can convey more information than words can de-
246 scribe. Users can obtain reliable thresholds for clustering since the effects of various thresholds
247 can be directly observed in the interface.

248 *4.2. Building a coordinate system*

249     To display a set of s-clusters in a space, we need to construct a coordinate system such that
250 interference among lines (different s-clusters) can be minimized in order to observe closeness.
251 The procedure is as follows. First, we examine the attribute value with the highest proportion
252 for each cluster. Then, summarize the number of distinct attribute values for each attribute,
253 and then sort them in increasing order. Attributes with the same number of distinct attribute val-
254 ues are further ordered by the lowest value of their proportions. The attributes with the least num-
255 ber of attribute values are arranged in the middle of the $X$-axis and others are put at two ends
256 according to the order described above. In other words, if the attribute values with the highest
257 proportion for all s-clusters are the same for some attribute $A_k$, this attribute will be arranged
258 in the middle of the $X$-axis. The next two attributes are then arranged at the left and right of $A_k$.
259     After the locations of attributes on the $X$-axis are decided, the locations of the corresponding
260 attribute values on the $Y$-axis are arranged accordingly. For each s-cluster, we examine the attri-
261 bute value with the highest proportion for each attribute. If the attribute value has not been seen
262 before, it is added to the "presenting list" (initially empty) for that attribute. Each attribute value
263 in the presenting list has a location as its order in the list. That is, not every attribute value has a
264 location on the $Y$-axis. Only attribute values with the highest proportion for some clusters have
265 corresponding locations on the $Y$-axis. Finally, we represent a s-cluster $C_x$ by its $d$ coordinates
266 $(L_x(i), L_y(v_{i,k}), P(v_{i,k}|C_x))$ for $i = 1, \ldots, d$, where the function $L_x(i)$ returns the $X$-coordinate for at-
267 tribute $A_i$, and $L_y(v_{i,k})$ returns the $Y$-coordinate for attribute value $v_{i,k}$.
268     In Fig. 4 for example, two s-clusters and their attribute distributions are shown in (a). Here, the
269 number of distinct attribute values with the highest proportion is 1 for all attributes except for $A_2$
270 and $A_7$. For these attributes, they are further ordered by their lowest proportions. Therefore, the
271 order for these eight attributes are $A_5, A_6, A_8, A_1, A_3, A_4, A_7, A_2$. With $A_5$ as center, $A_6$ and $A_8$ are
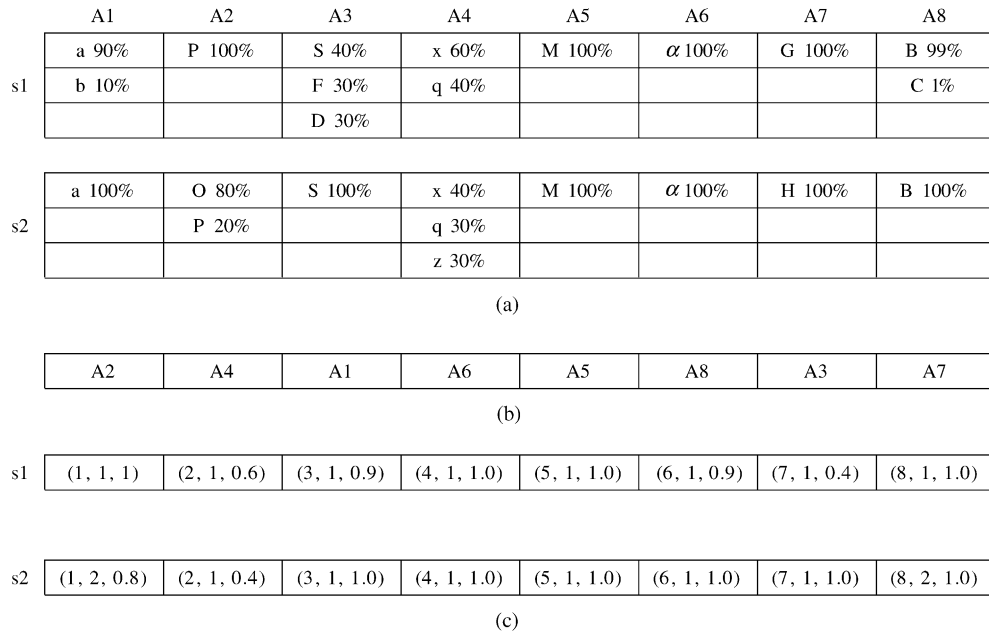
9

| | A1 | A2 | A3 | A4 | A5 | A6 | A7 | A8 |
|---|---|---|---|---|---|---|---|---|
| s1 | a 90% | P 100% | S 40% | x 60% | M 100% | $\alpha$ 100% | G 100% | B 99% |
| | b 10% | | F 30% | q 40% | | | | C 1% |
| | | | D 30% | | | | | |
| s2 | a 100% | O 80% | S 100% | x 40% | M 100% | $\alpha$ 100% | H 100% | B 100% |
| | | P 20% | | q 30% | | | | |
| | | | | z 30% | | | | |

(a)

| A2 | A4 | A1 | A6 | A5 | A8 | A3 | A7 |
|---|---|---|---|---|---|---|---|

(b)

| s1 | (1, 1, 1) | (2, 1, 0.6) | (3, 1, 0.9) | (4, 1, 1.0) | (5, 1, 1.0) | (6, 1, 0.9) | (7, 1, 0.4) | (8, 1, 1.0) |
|---|---|---|---|---|---|---|---|---|
| s2 | (1, 2, 0.8) | (2, 1, 0.4) | (3, 1, 1.0) | (4, 1, 1.0) | (5, 1, 1.0) | (6, 1, 1.0) | (7, 1, 1.0) | (8, 2, 1.0) |

(c)

Fig. 4. Example of constructing a coordinate system: (a) two s-clusters and their distribution table; (b) rearranged *X*-coordinate; (c) 3D-coordinates for s1 and s2.

arranged to the left and right, respectively. The rearranged order of attributes is shown in Fig. 4(b). Finally, we transform cluster $s_1$, and then $s_2$ into the coordinate system we build, as shown in Fig. 4(c). Taking $A_2$ for example, there are two attribute values *P* and *O* to be presented. Therefore, *P* gets a location 1 and *O* a location 2 at *Y*-axis. Similarly, *G* gets a location 1 and *H* a location 2 at *Y*-axis for $A_7$.

Fig. 5 shows an example of three s-clusters displayed in one window before (a) and after (b) the attribute rearrangement. The thicknesses of lines reflect the size of the s-clusters. Compared to the coordinate system without rearranging attributes, s-clusters are easier to observe in the new coordinate system since common points are located at the center along the *X*-axis presenting a trunk for the displayed s-clusters. For dissimilar s-clusters, there will be a small number of common points, leading to a short trunk. This is an indicator whether the displayed clusters are similar and this concept will be used in the interactive analysis described next.

### 4.3. Interactive visualization and analysis

The CDCS's interface, as described above, is designed to display the merging result of s-clusters such that users know the effects of adjusting merging parameters. However, instead of showing all s-clusters in the first step, our visualization tool displays only two groups from the merging result. More specifically, our visualization tool presents two groups in two windows for observing. The first window displays the group with the highest number of s-clusters since this group is usually the most complicated case. The second window displays the group which contains the cluster pair with the lowest similarity. The coordinate systems for the two groups are conducted respectively.

*C.-H. Chang, Z.-K. Ding / Data & Knowledge Engineering xxx (2004) xxx–xxx*
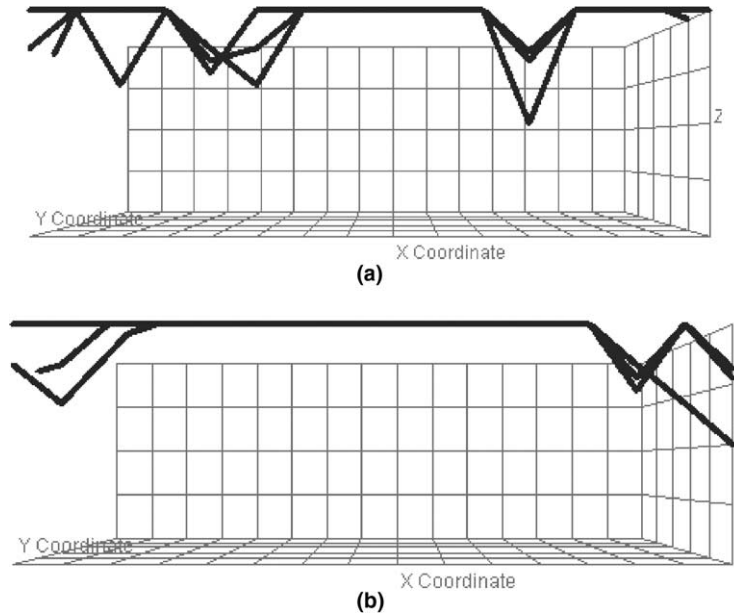


Fig. 5. Three s-clusters: (a) before and (b) after attribute rearrangement.

292   Fig. 6 shows an example of the CDCS's interface. The data set used is the Mushroom database
293 taken from the UCI machine learning repository [1]. The number of s-clusters obtained from the
294 first step is 106. The left window shows the group with the largest number of s-clusters, while the
295 right window shows the group with the least similar s-cluster pair. The number of s-clusters for
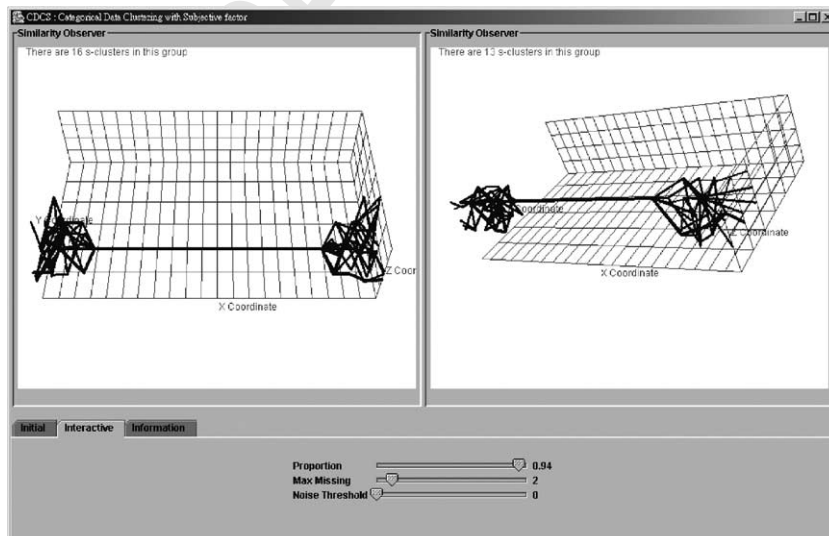


Fig. 6. Visualization of the mushroom data set ($e' = 2$).

296 these groups are 16 and 13, respectively, as shown at the top of the windows. Below these two
297 windows, three sliders are used to control the merging parameters for group merging and visual-
298 ization. The first two sliders denote the parameters $p'$ and $e'$ used to control the similarity thresh-
299 old $g'$. The third slider is used for noise control in the visualization so that small s-clusters can be
300 omitted to highlight the visualization of larger s-clusters. Each time the slider is moved, the binary
301 matrix BM is recomputed and merging result is updated in the windows. Users can also lock one
302 of the windows for comparison with a different threshold.
303     A typical process for interactive visualization analysis with CDCS is as follows. We start from a
304 strict threshold $g'$ such that the displayed groups are compact; and then relax the similarity thresh-
305 old until the displayed groups are too complex and the main trunk gets too short. A compact
306 group usually contains a long trunk such that all s-clusters in the group have the same values
307 and high proportions for these attributes. A complex group, on the other hand, presents a short
308 trunk and contains different values for many attributes. For example, both groups displayed in
309 Fig. 6 have obvious trunks which are composed of sixteen common points (or attribute values).
310 For a total of 22 attributes, 70% of the attributes have the same values and proportions for all s-
311 clusters in the group. Furthermore, the proportions of these attribute values are very high.
312 Through rotation, we also find that the highest proportion of the attributes on both sides of
313 the trunk is similarly low for all s-clusters. This implies that these attributes are not common fea-
314 tures for these s-clusters. Therefore, we could say both these groups are very compact since these
315 groups are composed of s-clusters that are very similar.
316     If we relax the parameter $e'$ from 2 to 5, the largest group and the group with least similar s-
317 clusters refer to the same group which contains 46 s-clusters, as shown in Fig. 7. For this merging
318 threshold, there is no obvious trunk for this group; and some of the highest proportions near the
319 trunk are relatively high, while others are relatively low. In other words, there are no common
320 features for these s-clusters, and thus this merge threshold is too relaxed since different s-clusters
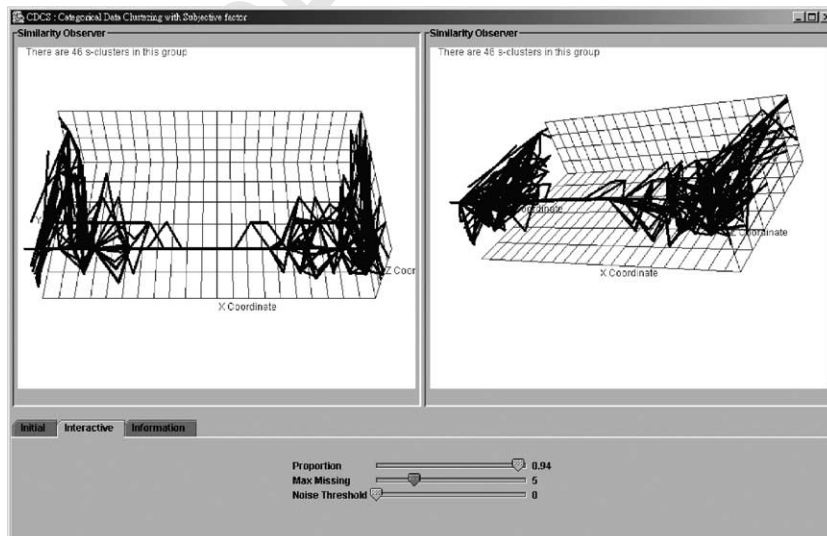


Fig. 7. Visualization of the mushroom data set with a mild threshold $e' = 5$.

12                  *C.-H. Chang, Z.-K. Ding / Data & Knowledge Engineering xxx (2004) xxx–xxx*
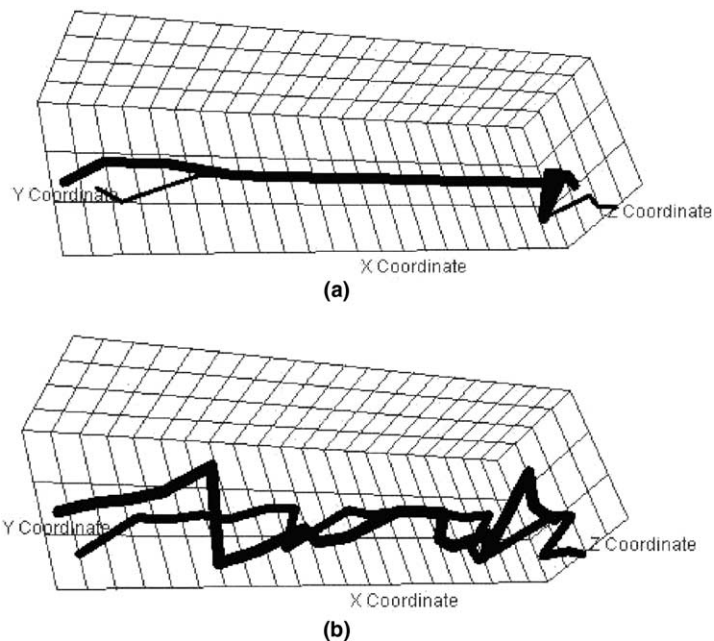


Fig. 8.  The most similar group pair for various merging threshold: (a) > (b).

321  are put in the same group. Therefore, the merging threshold in Fig. 6 is better than the one in
322  Fig. 7.

323   In summary, whether the s-clusters in a group are similar is based on users' viewpoints on the
324  obvious trunk. As the merging threshold is relaxed, more s-clusters are grouped together and the
325  trunks of both windows get shorter. Sometimes, we may reach a stage where the merge result is
326  the same no matter how the parameters are adjusted. This may be an indicator of a suitable clus-
327  tering result. However, it depends on how we view these clusters since there may be several such
328  stages. More discussion on this problem is presented in Section 5.2.

329   Omitting other merged groups does not do any harm since the smaller groups and the more
330  similar groups often have more common attribute values than the largest and the least similar
331  groups. However, to give a global view of the merged result, CDCS also offers a setting to display
332  all groups or a set of groups in a window. Particularly, we present the group pair with the largest
333  similarity as shown in Fig. 8, where (a) presents a high merging threshold and (b) shows a low
334  merging threshold for the data set mushroom. In principle, these two groups must disagree in
335  a certain degree or they might be merged by reducing the merging threshold. Therefore, users de-
336  cide the right parameter setting by finding the balance point where the displayed complex clusters
337  have long trunk, while the most similar group pair has very short or no trunk at all.

338  **5. Cluster validation**

339   Clustering is a field of research where its potential applications pose their own special require-
340  ments. Some typical requirements of clustering include scalability, minimal requirements for do-

341 main knowledge to determine input parameters, ability to deal with noisy data, insensitivity to the
342 order of input records, high dimensionality, interpretability and usability, etc. Therefore, it is
343 desirable that CDCS is examined under these requirements.

344
- First, in terms of scalability, the execution time of the CDCS algorithm is mainly spent on the first step. Simple cluster seeking requires only one database scan. Compared to EM-based algorithm such as AutoClass, which requires hundreds of iterations, this is especially desirable when processing large data sets.
- Second, the interactive visualization tool requires users with low domain knowledge to determine merging parameters.
- Third, the probability-based computation of similarity between objects and clusters can be easily extended to higher dimensions. Meanwhile, the clusters of CDCS can be simply described by the attribute–value pairs of high frequencies, suitable for conceptual interpretation.
- Finally, simple cluster seeking is sensitive to the order of input data, especially for skewed data. One way to alleviate this effect is to set a larger similarity threshold. The effect of parameter setting will be discussed in Section 5.3.

356
357 In addition to the requirements discussed above, the basic objective of clustering is to discover
358 significant groups present in a data set. In general, we should search for clusters whose members
359 are close to each other and well separated. The early work on categorical data clustering [9]
360 adopted an *external criterion* which measures the degree of correspondence between the clusters
361 obtained from our clustering algorithms and the classes assigned a priori. The proposed measure,
362 *clustering accuracy*, computes the ratio of correctly clustered instances of a clustering and is de-
363 fined as

$$\frac{\sum_{i=1}^{k} c_i}{n} \tag{5}$$

366 where $c_i$ is the largest number of instances with the same class label in cluster *i*, and *n* is the total
367 number of instances in the data set.
368 Clustering accuracy is only an indication of the intra-cluster consensus since high clustering
369 accuracy is easily achieved for larger numbers of clusters. Therefore, we also define two measures
370 using the data's interior criterion. First, we define intra-cluster cohesion for a clustering result as
371 the weighted cohesion of each cluster, where the cohesion for a cluster $C_k$ is the summation of the
372 highest probability for each dimension as shown below:

$$\text{intra} = \frac{\sum \text{in}_k * |C_k|}{n}, \quad \text{in}_k = \frac{\sum_{i=1}^{d} (\max_j P(v_{i,j}|C_k))^3}{d} \tag{6}$$

375 We also define inter-cluster similarity for a clustering result as the summation of cluster similarity
376 for all cluster pairs, weighted by the cluster size. The similarity between two clusters, $C_x$ and $C_y$, is
377 as defined in Eq. (3). The exponent $1/d$ is used for normalization since there are *d* component mul-
378 tiplications when computing $\text{Sim}(C_x, C_y)$.

$$\text{inter} = \frac{\sum_x \sum_y \text{Sim}(C_x, C_y)^{1/d} * |C_x \cup C_y|}{(k-1) * n} \tag{7}$$

14                    *C.-H. Chang, Z.-K. Ding / Data & Knowledge Engineering xxx (2004) xxx–xxx*

381    We present an experimental evaluation of CDCS on five real-life data sets from the UCI ma-
382 chine learning repository [1]. Four users are involved in the visualization analysis to decide a prop-
383 er grouping criterion. To study the effect due to the order of input data, each data set is randomly
384 ordered to create four test data sets for CDCS. The result is compared to AutoClass [2] and k-
385 mode [9], where the number of clusters required for k-mode is obtained from the clustering result
386 of CDCS.

387 *5.1. Clustering quality*

388    The five data sets used are Mushroom, Soybean-small, Soybean-large, Zoo and Congress vot-
389 ing, which have been used for other clustering algorithms. The size of the data set, the number of
390 attributes and the number of classes are described in the first column of Table 1. The mushroom
391 data set contains two class labels: poisonous and edible and each instance has 22 attributes. The
392 soybean-small and soybean-large contains 47 and 307 instances, respectively and each instance is
393 described by 35 attributes. (For soybean-small, there are 14 attributes which have only one value,
394 therefore these attributes are removed.) The number of class labels for soybean-small and soy-
395 bean-large are four and 19, respectively. The zoo data set contains 17 attributes for 101 animals.

Table 1
Number of clusters and clustering accuracy for three algorithms

| | # of clusters | | Accuracy | | |
|---|---|---|---|---|---|
| | AutoClass | CDCS | AutoClass | k-Mode | CDCS |
| Mushroom | 22 | 21 | 0.9990 | 0.9326 | 0.996 |
| 22 attributes | 18 | 23 | 0.9931 | 0.9475 | 1.0 |
| 8124 data | 17 | 23 | 0.9763 | 0.9429 | 0.996 |
| 2 labels | 19 | 22 | 0.9901 | 0.9468 | 0.996 |
| Zoo | 7 | 7 | 0.9306 | 0.8634 | 0.9306 |
| 16 attributes | 7 | 8 | 0.9306 | 0.8614 | 0.9306 |
| 101 data | 7 | 8 | 0.9306 | 0.8644 | 0.9306 |
| 7 labels | 7 | 9 | 0.9207 | 0.8832 | 0.9603 |
| Soybean-small | 5 | 6 | 1.0 | 0.9659 | 0.9787 |
| 21 attributes | 5 | 5 | 1.0 | 0.9361 | 0.9787 |
| 47 data | 4 | 5 | 1.0 | 0.9417 | 0.9574 |
| 4 labels | 6 | 7 | 1.0 | 0.9851 | 1.0 |
| Soybean-large | 15 | 24 | 0.664 | 0.6351 | 0.7500 |
| 35 attributes | 5 | 28 | 0.361 | 0.6983 | 0.7480 |
| 307 data | 5 | 23 | 0.3224 | 0.6716 | 0.7335 |
| 19 labels | 5 | 21 | 0.3876 | 0.6433 | 0.7325 |
| Congress voting | 5 | 24 | 0.8965 | 0.9260 | 0.9858 |
| 16 attributes | 5 | 28 | 0.8942 | 0.9255 | 0.9937 |
| 435 data | 5 | 26 | 0.8804 | 0.9312 | 0.9860 |
| 2 labels | 5 | 26 | 0.9034 | 0.9308 | 0.9364 |
| Average | | | 0.8490 | 0.8716 | 0.9260 |

Table 2
Comparison of AutoClass and CDCS

| Data set | Intra | | Inter | | Intra/inter | |
|---|---|---|---|---|---|---|
| | AutoClass | CDCS | AutoClass | CDCS | AutoClass | CDCS |
| Mushroom | 0.6595 | 0.6804 | 0.0352 | 0.0334 | 18.7304 | *20.3704 |
| Zoo | 0.8080 | 0.8100 | 0.1896 | 0.2073 | *4.2663 | 3.9070 |
| Soybean-small | 0.6593 | 0.7140 | 0.1840 | 0.1990 | 3.5831 | 3.5879 |
| Soybean-large | 0.5826 | 0.7032 | 0.1667 | 0.1812 | 3.4940 | *3.8807 |
| Congress voting | 0.5466 | 0.6690 | 0.1480 | 0.3001 | *3.6932 | 2.2292 |

396 After data cleaning, there are 16 attributes and each data object belonging to one of seven classes.
397 The Congress voting data set contains the votes of 435 congressman on 16 issues. The congress-
398 man are labelled as either Republican or Democratic.
399   Table 1 records the number of clusters and the clustering accuracy for the five data sets. As
400 shown in the last row, CDCS has better clustering accuracy than the other two algorithms. Fur-
401 thermore, CDCS is better than k-mode in each experiment given the same number of clusters.
402 Compared with AutoClass, CDCS has even more clustering accuracy since it finds more clusters
403 than AutoClass, especially for the last two data sets. The main reason for this phenomenon is that
404 CDCS reflects the user's view on the degree of intra-cluster cohesion. Various clustering results,
405 say nine clusters and 10 clusters, are not easily observed in this visualization method. Therefore,
406 if we look into the clusters generated by these two algorithms, CDCS has better intra-cluster cohe-
407 sion for all data sets; whereas AutoClass has better cluster separation (smaller inter-cluster sim-
408 ilarity) on the whole, as shown in Table 2. In terms of intra-cluster similarity over inter-cluster
409 similarity, AutoClass performs better on Zoo and Congress voting, whereas CDCS performs bet-
410 ter on Mushroom and Soybean-large.

411 *5.2. Discussion on cluster numbers*

412   To analyze the data sets further, we record the number of clusters for each merging threshold of
413 the second step. The merging thresholds, $g'$, are calculated by Eq. (4), where $p'$ and $e'$ vary from 0
414 to 0.99 and 1 to 4, respectively. A total of 400 merging thresholds are sorted in decreasing order.
415 The number of clusters for each merging threshold is recorded until the number clusters reaches
416 three. The way the merging thresholds are calculated avoids steep curves where the number of
417 clusters changes rapidly with small merging thresholds. As shown in Fig. 9(a), we can see five
418 smooth curves with steep downward slopes at the zero ends. The small merging thresholds are
419 a result of the similarity function between two s-clusters (Eq. (3)) where a series multiplications
420 are involved (each factor represents the percentage of common values of an attribute). Therefore,
421 we also change the scale of the merging threshold $g'$ to $\sqrt[d]{g'}$ ($d$ is the number of dimensions), which
422 represents the average similar of an attribute as shown in Fig. 9(b).
423   We try to seek smooth fragments for each curve where the number of clusters does not change
424 with the varying merging thresholds. Intuitively, these smooth levels may correspond to some
425 macroscopic views where the number of clusters are persuasive. For example, the curve of Zoo
426 has a smooth level when the number of clusters are 9, 8, 7, etc., Mushroom has a smooth level
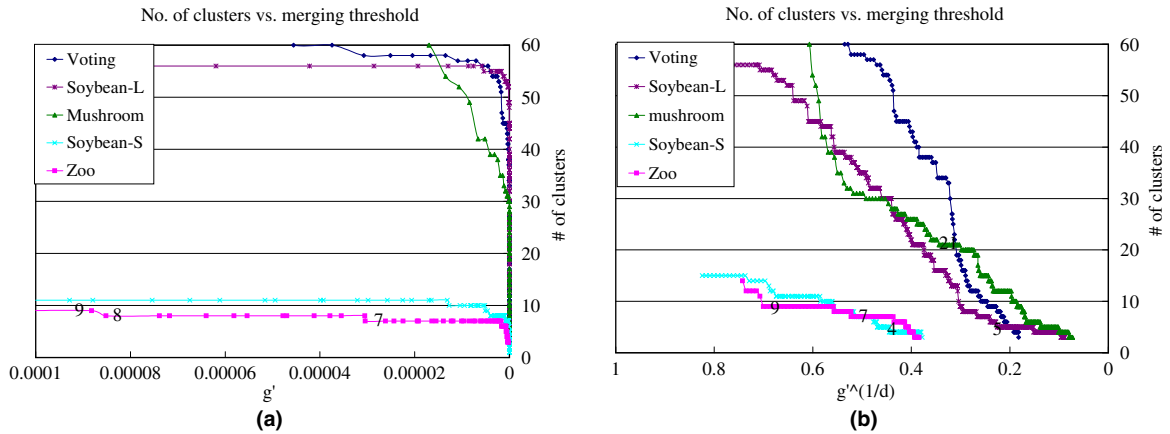
Fig. 9. Number of clusters vs. merging threshold.

427 at 30, 21, 20, and 12, while the curve of Soybean-small has a smooth level at 7, 5, 4, etc. Some of
428 these coincide with the clustering results of AutoClass and CDCS. Note that the longest fragment
429 does not necessarily symbolize the best clustering result since it depends on how we compute the
430 similarity and the scale of the *X*-axis.

431   We find that Soybean-large has a smooth fragment when the number of cluster is 5 which cor-
432 responds to that of AutoClass, however the average similarity of an attribute is dropped below
433 0.22. We also notice that the Congress voting has quite steep slope at cluster number between
434 30 and 5. These may indicate that the data set itself contains a complex data distribution such that
435 no obvious cluster structure is present. We believe that the optimal cluster number varies with dif-
436 ferent clustering criteria, a decision for the users. For Congress voting, high clustering accuracy is
437 more easily achieved since the number of class labels is only two. As for Soybean-large, clustering
438 accuracy cannot be high if the number of clusters is less than 19. Note that class labels are given
439 by individuals who categorize data based on background domain knowledge. Therefore, some at-
440 tributes may weight more heavily than others. However, most computations of the similarity be-
441 tween two objects or clusters gives equal weight to each attribute. We intend to extend our
442 approach to investigate these further in the future.

443 *5.3. Parameter setting for dynamic clustering*

444   In this section, we report experiments on the first step to study the effect of input order and
445 skewed cluster size versus various threshold setting. For each data set, we prepare 100 random
446 orders of the same data set and run the dynamic clustering algorithms 100 times for $p = 0.9$,
447 $p = 0.8$, $p = 0.7$, respectively. The mean number of clusters as well as the standard deviation are
448 shown in Table 3. The number of clusters generated increases as the merging threshold $p$ increases.
449 For mushroom, the mean number of clusters varies from 49 to 857 when $p$ varies from 0.7 to 0.9.
450   To study the effect of skewed cluster sizes, we conduct the following experiment: we run the sim-
451 ple clustering twice with the data order reversed for the second run. To see if the clustering is more
452 or less the same, we use a confusion table $V$ with as rows the clusters in the first run, and as col-

Table 3
The mean number of clusters and its standard deviation for various *p*

| # of clusters | *p* = 0.7 | | *p* = 0.8 | | *p* = 0.9 | |
|---|---|---|---|---|---|---|
| | Mean | S.D. | Mean | S.D. | Mean | S.D. |
| Mushroom | 49.09 | 2.29 | 165.75 | 4.16 | 857.98 | 6.75 |
| Zoo | 9.34 | 1.32 | 12.96 | 2.06 | 16.28 | 2.47 |
| Soybean-small | 16.14 | 1.67 | 21.66 | 2.46 | 26.296 | 2.87 |
| Soybean-large | 77.64 | 2.29 | 108.71 | 5.66 | 151.26 | 8.11 |
| Congress voting | 69.62 | 2.33 | 97.27 | 5.33 | 130.90 | 7.30 |

453 umns the clusters in the second run. Then, the entry *i, j* corresponds to the number of data objects
454 that were in cluster *i* in the first experiment, and in cluster *j* in the second experiment.

455    The consistency of the two runs can be measured by the percentage of zero entries in the con-
456 fusion table. However, since two large numbers of small clusters tend to have large zero percent-
457 age compared to two small numbers of large clusters, these values are further normalized by the
458 largest number of possible values. For example, consider an $m \times n$ confusion table, where $m$ and $n$
459 denotes the number of clusters generated for the two runs of the reverse order experiment. The
460 largest number of zero entries will be $(m * n - \max\{m, n\})$.

461    Let $V$ denotes the confusion table for the two runs. Therefore, the normalized zero percentage
462 (NZP) of the confusion table is defined by

$$\text{NZP}(V) = \frac{|\{(i, j)|V(i, j) = 0, 1 \leqslant i \leqslant m, 1 \leqslant j \leqslant n\}|}{m * n - \max\{m, n\}} \tag{8}$$

465 We show the mean NZP and its standard deviation of 100 reverse order experiments for $p = 0.7$
466 and $p = 0.9$ with the five data sets in Table 4. For each data set, we also prepare a skewed input
467 order where the data objects of the largest class are arranged one by one, then come the data ob-
468 jects of the second largest class, and so on. The NZP and the numbers of clusters for the reverse

Table 4
Comparison of NZP for skew input order and random order

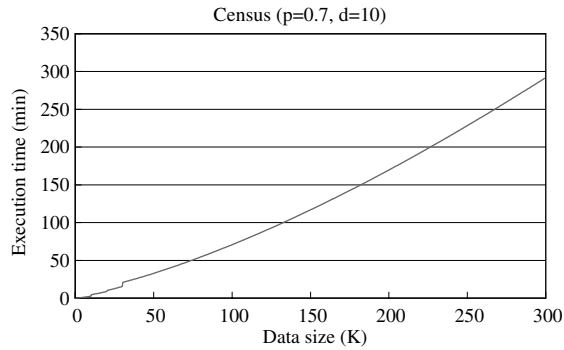| NZP | Rand vs. Rand | | Skew vs. Skew | | Rand vs. Skew | |
|---|---|---|---|---|---|---|
| | Mean | S.D. | NZP | $m \times n$ | Mean | S.D. |
| *(a) p = 0.7* | | | | | | |
| Mushroom | 0.9904 | 0.0372 | 0.9962 | 44 × 25 | 0.9894 | 0.0244 |
| Zoo | 0.9036 | 0.1676 | 0.9166 | 10 × 7 | 0.9100 | 0.1229 |
| Soybean-small | 0.9658 | 0.0927 | 0.9553 | 16 × 15 | 0.9665 | 0.0917 |
| Soybean-large | 0.9954 | 0.0282 | 0.9949 | 85 × 66 | 0.9945 | 0.0244 |
| Congress voting | 0.9887 | 0.0398 | 0.9878 | 69 × 50 | 0.9904 | 0.0281 |
| *(b) p = 0.9* | | | | | | |
| Mushroom | 0.9986 | 0.0071 | 0.9980 | 489 × 423 | 0.9986 | 0.0100 |
| Zoo | 0.9857 | 0.0695 | 0.9893 | 26 × 19 | 0.9869 | 0.0620 |
| Soybean-small | 0.9944 | 0.0360 | 0.9956 | 33 × 29 | 0.9963 | 0.0373 |
| Soybean-large | 0.9993 | 0.0283 | 0.9994 | 233 × 225 | 0.9994 | 0.0077 |
| Congress voting | 0.9986 | 0.0100 | 0.9986 | 196 × 161 | 0.9987 | 0.0100 |

Fig. 10.  The execution time required for simple cluster seeking.

469 order experiments of the skewed input order are displayed in the middle columns (Skew vs. Skew
470 column) of the table for comparison with the average cases (Rand vs. Rand column) for each data
471 set. We also show the mean NZP and its variance between the skew input order and each random
472 order (Rand vs. Skew column). From the statistics, there is no big difference in the skewed input
473 order and average cases. Comparing various $p$: (a) 0.7 and (b) 0.9, the mean NZP increases as the
474 parameter $p$ increases. This validates our claim in Section 3 that the effect of input order decreases
475 as the threshold increases.

476     Finally, we use the Census data set in the UCI KDD repository for scalability experiment. This
477 data set contains weighted census data extracted from the 1994 and 1995 current population sur-
478 veys conducted by the US Census Bureau. There are a total of 199,523 + 99,762 instances, each
479 with 41 attributes. Fig. 10 shows the execution time for simple cluster seeking of the Census data
480 set with increasing data size. It requires a total of 300 minutes to cluster 299,402 objects with
481 parameter setting $p = 0.7$ and $e = 10$. Note that CDCS is implemented using Java and no code
482 optimization has been used. The clustering accuracy for the first step is 0.9406. For the cluster
483 merging step, it costs 260 seconds to merge 5865 s-clusters for the first visualization, and 64 sec-
484 onds for the second visualization. After users' subjective judgement, a total of 359 clusters are gen-
485 erated. Comparing to AutoClass, it takes more than two days before it completes the clustering.

486 **6. Conclusion and future work**

487     In this paper, we introduced a novel approach for clustering categorical data with visualization
488 support. First, a probability-based concept is incorporated in the computation of object similarity
489 to clusters; and second, a visualization method is devised for presenting categorical data in a 3D
490 space. Through an interactive visualization interface, users can easily decide a proper parameter
491 setting. Thus, human subjective adjustment can be incorporated in the clustering process. From
492 the experiments, we conclude that CDCS performs quite well compared to state-of-the-art clus-
493 tering algorithms. Meanwhile, CDCS successfully handles data sets with significant differences
494 in the sizes of clusters such as Mushroom. In addition, the adoption of naive-Bayes classification
495 makes CDCS's clustering results much more easily interpreted for conceptual clustering.

This visualization mechanism may be adopted for other clustering algorithms which require parameter adjustment. For example, if the first step is replaced by complete-link hierarchical clustering with high similarity threshold, we will be able to apply the second step and the visualization technique to display the clustering result and let users decide a proper parameter setting. Another feature that can be included in CDCS is the figure of some clustering validation measure versus our merging threshold. Such measures will enhance users' confidence on the clustering result. In the future, we intend to devise another method to enhance the visualization of different clusters. Also, we will improve the CDCS algorithm to handle data with both categorical and numeric attributes.

## Acknowledgement

## References

[1] C.L. Blake, C.J. Merz, UCI repository of machine learning databases, <http://www.cs.uci.edu/~mlearn/MLRepository.html>, Department of Information and Computer Science, University of California, Irvine, CA, 1998.

[2] P. Cheeseman, J. Stutz, Bayesian classification (autoclass): theory and results, in: Proceedings of Advances in Knowledge Discovery and Data Mining, 1996, pp. 153–180.

[3] D. Fisher, Improving inference through conceptual clustering, in: Proceedings of AAAI-87 Sixth National Conference on Artificial Intelligence, 1987, pp. 461–465.

[4] M. Friendly. Visualizing categorical data: data, stories, and pictures, in: SAS Users Group International, 25th Annual Conference, 2002.

[5] V. Ganti, J. Gehrke, R. Ramakrishnan, Cactus—clustering categorical data using summaries, in: Proceedings of the fifth ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, 1999, pp. 73–83.

[6] D. Gibson, J. Kleinberg, P. Raghavan, Clustering categorical data: an approach based on dynamical systems, VLDB Journal 8 (1998) 222–236.

[7] S. Guha, R. Rastogi, K. Shim, Rock: a robust clustering algorithm for categorical attributes, Information Systems 25 (2000) 345–366.

[8] E.-H. Han, G. Karypis, V. Kumar, B. Mobasher, Clustering based on association rule hypergraphs, in: Workshop on Research Issues in Data Mining and Knowledge Discovery (DMKD), 1997, pp. 343–348.

[9] Z. Huang, Extensions to the k-means algorithm for clustering large data sets with categorical values, Data Mining and Knowledge Discovery 2 (1998) 283–304.

[10] T. Kohonen, Self-organizing Maps, Springer-Verlag, 1995.

[11] T. Kohonen, S. Kaski, K. Lagus, T. Honkela, Very large two-level SOM for the browsing of newsgroups, in: Proceedings of International Conference on Artificial Neural Networks (ICANN), 1996, pp. 269–274.

[12] A. Konig, Interactive visualization and analysis of hierarchical neural projections for data mining, IEEE Transactions on Neural Networks 11 (3) (2000) 615–624.

[13] W.A. Kosters, E. Marchiori, A.A.J. Oerlemans, Mining clusters with association rules, in: Proceedings of Advances in Intelligent Data Analysis, 1999, pp. 39–50.

[14] S. Ma, J.L. Hellstein, Ordering categorical data to improve visualization, in: IEEE Symposium on Information Visualization, 1999.

537  [15] J.B. MacQueen, Some methods for classification and analysis of multivariate observations, in: Proceedings of the
538       5th Berkeley Symposium on Mathematical Statistics and Probability, pp. 281–297, 1967.
539  [16] P.C. Mahalanobis, Proceedings of the National Institute of Science of India 2 (49) (1936).
540  [17] T.M. Mitchell, Machine Learning, McGraw-Hill, 1997.
541  [18] C.J. van Rijsbergen, Information Retrieval, Butterworths, London, 1979 (Chapter 3).
542  [19] E. Sirin, F. Yaman, Visualizing dynamic hierarchies in treemaps. <http://www.cs.umd.edu/class/spring2002/
543       cmsc838f/Project/DynamicTreemap.pdf>, 2002.
544  [20] J.T. To, R.C. Gonzalez, Pattern Recognition Principles, Addison-Wesley Publishing Company, 1974.
545  [21] A.K.H. Tung, J. Hou, J. Han, Spatial clustering in the presence of obstacles, in: Proceedings of 2001 International
546       Conference on Data Engineering, 2001, pp. 359–367.
547  [22] K. Wang, C. Xu, B. Liu, Clustering transactions using large items, in: Proceedings of the ACM CIKM
548       International Conference on Information and Knowledge Management, 1999, pp. 483–490.
549  [23] Y. Zhang, A. Fu, C.H. Cai, P. Heng, Clustering categorical data, in: Proceedings of 16th IEEE International
550       Conference on Data Engineering, 2000, p. 305.
551

**Chia-Hui Chang** is an assistant professor at the Department of Computer Science and Information Engineering, National Central University in Taiwan. She received her B.S. in Computer Science and Information Engineering from National Taiwan University, Taiwan in 1993 and got her Ph.D. in the same department in January 1999. Her research interests include information extraction, data mining, machine learning, and Web related research. Her URL is http://www.csie.ncu.edu.tw/~chia/.

**Zhi-Kai Ding** received the B.S. in Computer Science and Informantion Engineering from National Dong-Hwa University, Taiwan in 2001, and M.S. in Computer Science and Informantion Engineering from National Central University, Taiwan in 2003. Currently, he is working as a software engineer at Hyweb Technology Co., Ltd. in Taiwan. His research interest includes data mining, information retrieval and extraction.