

Research Article

A Novel Method of Interestingness Measures for Association Rules Mining Based on Profit

Chunhua Ju,^{1,2} Fuguang Bao,^{1,3} Chonghuan Xu,³ and Xiaokang Fu^{1,3}

¹Contemporary Business and Trade Research Center, Zhejiang Gongshang University, Hangzhou 310018, China

²College of Computer Science & Information Engineering, Zhejiang Gongshang University, Hangzhou 310018, China

³School of Business Administration, Zhejiang Gongshang University, Hangzhou 310018, China

Correspondence should be addressed to Fuguang Bao; baofuguang@126.com

Received 8 April 2015; Revised 17 June 2015; Accepted 18 June 2015

Academic Editor: Allan C. Peterson

Copyright © 2015 Chunhua Ju et al. This is an open access article distributed under the Creative Commons Attribution License, which permits unrestricted use, distribution, and reproduction in any medium, provided the original work is properly cited.

Association rules mining is an important topic in the domain of data mining and knowledge discovering. Some papers have presented several interestingness measure methods; the most typical are *Support*, *Confidence*, *Lift*, *Improve*, and so forth. But their limitations are obvious, like no objective criterion, lack of statistical base, disability of defining negative relationship, and so forth. This paper proposes three new methods, *Bi-lift*, *Bi-improve*, and *Bi-confidence*, for *Lift*, *Improve*, and *Confidence*, respectively. Then, on the basis of utility function and the executing cost of rules, we propose interestingness function based on profit (*IFBP*) considering subjective preferences and characteristics of specific application object. Finally, a novel measure framework is proposed to improve the traditional one through experimental analysis. In conclusion, the new methods and measure framework are prior to the traditional ones in the aspects of objective criterion, comprehensive definition, and practical application.

1. Introduction

Along with the development and progress of data collection and storage technology in economy society, various industries and fields gradually accumulate a lot of data. In order to extract valuable information from these data, researchers have paid more and more attention to data mining technology. Association rules mining is a key technology in data mining field and is applied in many areas [1]. Association rules algorithm can generate a lot of rules, but, due to limited resources, only a part of the rules could be used by makers [2]. Therefore, the interestingness evaluation is significant for the practical application of association rules mining technology.

In recent years there have been quite a few results on interestingness measure of association rules [3, 4]. Tseng et al. put forward incremental maintenance of generalized association rules under taxonomy evolution [5, 6]. Literature [7, 8] separately use two different data envelopment analysis models for the evaluation of an association rules set. To compare the evaluation results, we can find that use of different evaluation methods will get different results from the same association rules set. Interestingness measures

mainly include objective measure and subjective measure [9]. The objective interestingness mainly focuses on the statistics significant research on objective data, the classic ones include *Support*, *Confidence*, and *Lift*, and the new one includes *Validity*, *Conviction*, and *Improve*. Subjective interestingness involves the personality characteristics of subject (users) such as domain knowledge and the hobbies. There is less research on subjective interestingness and it is relatively immature. Hoque et al. and Zhang et al. presented methods to generate both frequent and rare itemsets using multiobjective genetic algorithm [10, 11]. How to mine the real and effective association rules reflecting interests of users is the common goal for researchers.

The following defects of the association rules mining methods exist: (1) Many traditional association methods generate a lot of rules, and most of them are not relevant or even rules of error. (2) Do not consider users characteristics and their changes, but users characteristics and subjectivity tend to affect relevance of several events. (3) Online trading data and user evaluation data are extremely sparse (i.e., data sparseness) due to surge of current online trading. (4) The too low threshold values of *support* and *confidence* can produce

the combination explosion, but, because of data sparseness, low support rules may provide some novel knowledge that users are interested in. (5) At present some literature simply combines kinds of interestingness evaluation to measure, but this does not take rationality of various interestingness evaluation methods into account.

The main contribution of this paper is summarized as follows: This paper firstly summarizes the relevant research results on the objective interestingness measure of association rules. At the same time the various measurement methods of association rules are analyzed and compared, and we found that all have some defects and problems. And, then, put forward three more effective measurement methods of association rules (*Bi-lift*, *Bi-improve*, and *Bi-confidence*) based on improvement of some methods. The *Bi-lift* method takes deduction of negative premise as a constraint, to form a bideduction comparing algorithm so as to improve the reliability of the mutual influence between premise and follow-up. The *Bi-improve* method takes the adjustment of the occurrence and nonoccurrence possibility of antecedent based on the *Improve*. The *Bi-confidence* method takes the adjustment of the nonoccurrence possibility of antecedent based on the *Confidence*. Through the experimental analysis, a new measure framework is proposed to improve the traditional one. Then, on the basis of utility function and the executing cost of rules, we propose interestingness function based on profit (*IFBP*) considering subjective preferences and characteristics of specific application object.

A review of related work is given in Section 2. In Section 3, three more effective measurement methods of association rules (*Bi-lift*, *Bi-improve*, and *Bi-confidence*) are proposed. In Section 4, the paper studied subjective interestingness measures of association rules. In Section 5, through the experimental analysis, a new measure framework is proposed to improve the traditional one. Finally, we draw conclusions in Section 6.

2. Review on Objective Measures of Interest for Association Rules

To illustrate conveniently, firstly we suppose that formal description of association rules is as follows:

$$A \longrightarrow B. \quad (1)$$

In this description, $A = \{A_1, A_2, \dots, A_j\} \subset I$ and $B = \{B_1, B_2, \dots, B_k\} \subset I$. I indicate itemsets, and $A \cap B = \phi$. Rules should meet certain *support* threshold s and *confidence* c .

Data of Table 1 is produced by data extraction and transformation from a shopping mall receipts. Each line (a tuple) is a shopping list data (shopping receipts). "1" represents that the list includes this item, while "0" represents that the list does not include this item.

2.1. Support. *Support* [12] means the frequency that the data fields A and B involved in association rules occur together in the data set. Only the association rules appear frequently in the itemsets, when it gets high accuracy. *Support* can be used to measure the usefulness of association rules. When the

TABLE 1: A data set of transaction.

Tid	E	F	G	H	I	J	K	L	M	N	R	Ttotal
1	1	1	1	1	1	1	1	0	1	0	0	8
2	1	1	0	1	1	1	1	1	1	0	0	8
3	1	1	0	1	1	0	0	0	0	1	0	5
4	1	1	1	0	0	1	0	1	1	1	1	8
5	1	1	1	0	0	1	0	0	0	1	1	6
6	1	1	1	0	0	0	0	0	1	1	1	6
7	1	1	0	0	0	0	1	0	0	1	0	4
8	1	1	1	0	1	0	0	0	0	0	1	5
9	1	0	1	1	0	0	0	0	0	1	0	4
10	1	0	1	0	0	1	0	0	0	0	1	4
Total	10	8	7	4	4	5	3	2	4	6	5	58

frequency of A and B occurring at the same time is equal to or greater than the designated minimum *support* threshold, A and B meet frequent itemsets. *Support* can be expressed as

$$s(A \longrightarrow B) = P(AB) = \frac{N(AB)}{|D|}, \quad (2)$$

where $N(AB)$ is the record number of A and B that appeared together, and $|D|$ is the total record number of transactions in data sets.

Support is classic but also has the defects of artificially controlled threshold and rare itemsets. Many infrequent itemsets in the data set may have potential value. Besides, at present in large electronic commerce system, the number of subjects (users) and the amount of projects increase exponentially. Online transaction data and user evaluation data are extremely sparse.

2.2. Confidence. *Confidence* [12] is the statistics of probability $P(B | A)$ that subsequent events occur under the condition of occurrence of the precursor events in trading data sets. It is used to measure the reliability of the rules. Formula is

$$c(A \longrightarrow B) = P(B | A) = \frac{P(AB)}{P(A)}. \quad (3)$$

It is used to combine *confidence* with *support* to form *Support-confidence* framework for mining association rules [13]. If *Support* is larger than the designated minimum *support* threshold and *Confidence* is larger than the designated minimum *confidence* threshold, the rules are called strong association rules. But strong association rules are not always effective, some are not what users are interested in, and some are even misleading.

2.3. Lift. Because of the defects of *Support-confidence* framework, some scholars analyze the relativity of association rules mined, namely, *lift* [12]. *Lift* means the ratio of rule's *Confidence* to probability of occurrence of the consequent, which reflects positive or negative correlation of antecedent and consequence of rules. It refers to the ratio of the occurrence probability of B under the condition A to that without

TABLE 2: The occurrence of “A” and “B.”

	B occurring	B not occurring	Total
A occurring	4500	2000	6500
A not occurring	3000	500	3500
Total	7500	2500	10000

considering condition A, which reflects the relationship between “A” and “B”:

$$\text{lift}(A \rightarrow B) = \frac{c(A \rightarrow B)}{P(B)} = \frac{P(AB)}{P(A)P(B)}. \quad (4)$$

The range of *lift* values is $[0, +\infty)$. As *lift* is equal to 1, it shows that A and B appearing at the same time belong to independent random events and have no special significance; namely, A and B are independent of each other with no mutual affection. We call this rule uncorrelated rules; if *lift* value is less than 1, it shows that the emergence of “A” reduces the emergence of “B,” and then we call them negative correlation rules; if *Lift* value is larger than 1, it shows that the emergence of “A” promotes the emergence of “B,” and then we call them positive correlation rules. Problems: *Lift* takes events A and B in equivalence position. According to the *Lift*, $A \rightarrow B$ and $B \rightarrow A$ are the same; that is to say, if we accept rule $A \rightarrow B$, $B \rightarrow A$ should be also accepted, but the fact is not like this.

2.4. *Validity*. Literature [14] introduces a new measure method of association rules, known as *validity*. *Validity* is defined as the difference between the probability of “A” and “B” occurring together and the occurrence probability of “B” without “A” occurring in database D. Because the value range of $P(AB)$ and $P(\overline{AB})$ are $[0, 1]$, the value range of *validity* is obviously $[-1, 1]$:

$$\text{Validity}(A \rightarrow B) = P(AB) - P(\overline{AB}). \quad (5)$$

In fact, our research found that the *validity* is not effective. In Table 1, for example, rule $F \rightarrow G$'s *support* is 0.7 and its *validity* is $0.5 - 0.2 = 0.3$; according to the literature [14], we can judge that it is a valid association rule. But through the calculation $P(FG) - P(F)P(G) = 0.5 - 0.8 \times 0.7 = -0.06$, it shows that F and G have certain negative correlation. And taking Table 2 for example, this case has already qualified for the basic requirements of *support* and *confidence*, and $\text{Validity}(A \rightarrow B) = 0.45 - 0.3 = 0.15$. According to the measure standard of *validity*, the rule $A \rightarrow B$ should be said “effective.” The occurrence of “A” will promote the occurrence of “B”. But in fact, the occurrence frequency of overall event B is 0.75, and the occurrence possibility of “B” as “A” occurs is $4500/6500 = 0.69$, which is six percent lower.

2.5. *Conviction*. As early as 1997, Brin introduced the concept of *conviction* (*Conv.*) [15, 16]:

$$\text{Conviction}(A \rightarrow B) = \frac{P(A)P(\overline{B})}{P(\overline{AB})}. \quad (6)$$

TABLE 3: The occurrence of “A’” and “B’.”

	B’ occurring	B’ not occurring	Total
A’ occurring	8000	1000	9000
A’ not occurring	500	500	1000
Total	8500	1500	10000

TABLE 4: The occurrence of “C” and “D.”

	D occurring	D not occurring	Total
C occurring	3600	700	4300
C not occurring	3700	2000	5700
Total	7300	2700	10000

Its value range is $[0, +\infty)$. When the value of *conviction* is “1,” it means that “A” has no relation with “B.” And the greater the *conviction* is, the higher the interest in the rule will be. But the *conviction* constraints are too strict; lots of valuable association rules will be removed. In Table 1, for example, $\text{Conviction}(M \rightarrow L) = 0.4 \times 0.8 / 0.5 = 0.64 < 1$; its *conviction* is low, but, in fact, the high interest of M and L may exist. And it requires that $P(A\overline{B}) \neq 0$ at the same time.

2.6. *Improve*. Literature [17] proposed a new interestingness measure method of association rules based on the description of the defects of the traditional interestingness measurement method. We will call it “*Improve*.” It means that the difference of the conditional probability $P(B | A)$ and the probability of “B”

$$\text{Improve}(A \rightarrow B) = [P(B | A) - P(B)]. \quad (7)$$

But shortcomings of *Improve* (*Imp.*) are obvious. Firstly, how much improvement of probability can be called improvement? Secondly, the probability of former pieces’ occurrence will seriously affect *Improve* evaluation in such a way that when it is high, the *improve* value will be very small all the time.

Take Tables 3 and 4 for example and calculate their *Improve* values:

$$\text{Improve}(A' \rightarrow B') = [P(B' | A') - P(B')] = 0.03, \quad (8)$$

$$\text{Improve}(C \rightarrow D) = [P(D | C) - P(D)] = 0.11.$$

Only, thinking about the *Improve*, rule $(C \rightarrow D)$ is more valuable than rule $(A' \rightarrow B')$. But the fact is also very clear that it can increase “the occurrence possibility of B” as “A” occurs by up to 39% compared to that when “A” does not occur. While it can increase “the occurrence possibility of event D” as “C” occurs by up to 19% compared to that when “C” does not occur. So rule $(A' \rightarrow B')$ should be more meaningful than rule $(C \rightarrow D)$.

2.7. *Chi-Square Analysis*. Literature [18] puts forward an interestingness measure standard based on *t*-test. It uses *t*-test to analyze difference of associated *confidence* $P(B | A)$ and expected *confidence* $P(B)$. If the difference is bigger, it

indicates that the occurrence of “A” has large influence on “B,” and the rule $(A \rightarrow B)$ is interesting. Formula is as follows:

$$\begin{aligned} \text{Chi-Square}(A \rightarrow B) &= \frac{[P(B|A) - P(B)]}{\sigma}, \\ \sigma &= \sqrt{\frac{P(B)(1 - P(B))}{n}}. \end{aligned} \quad (9)$$

If $\text{Chi-Square}(A \rightarrow B) > t_{\alpha}(n)$, it shows that it has larger difference between associated *confidence* $P(B|A)$ and expected *confidence* $P(B)$ and the rule $(A \rightarrow B)$ is interesting. To some extent, it improves the traditional framework of interestingness measure. But there is also a defect that is similar to that of *Improve*.

2.8. *Certainty Factor (CF)*. To assess the accuracy of association rules, Berzal et al. use certainty factors [19] instead of *confidence*. Formula is as follows:

$$\begin{aligned} \text{CF}(A \rightarrow B) &= \frac{c(A \rightarrow B) - s(B)}{1 - s(B)} \\ &= \frac{P(AB)/P(A) - P(B)}{1 - P(B)} \\ &= \frac{P(AB) - P(A)P(B)}{P(A)(1 - P(B))}, \\ &\quad \text{if } c(A \rightarrow B) > s(B), \\ \text{CF}(A \rightarrow B) &= \frac{c(A \rightarrow B) - s(B)}{s(B)} \\ &= \frac{P(AB)/P(A) - P(B)}{P(B)} \\ &= \frac{P(AB) - P(A)P(B)}{P(A)P(B)}, \\ &\quad \text{if } c(A \rightarrow B) < s(B). \end{aligned} \quad (10)$$

The certainty factor is interpreted as a measure of variation of the probability that B is in a transaction when we consider only those transactions where there is A . More specifically, a positive *CF* measures the decrease of the probability that B is not in a transaction, given that A is. A similar interpretation can be done for negative *CF*.

3. The Improvement of Objective Interestingness Measures

3.1. *Bi-Lift*. Related researches show that *lift* method has good evaluation results. But obviously *lift* take “A” and “B” in equivalence position, and it shows that rules $A \rightarrow B$ and $B \rightarrow A$ are the same; if we accept rule $A \rightarrow B$, rule $B \rightarrow A$ should be also accepted. But the fact is not like this. For this problem, the paper proposes a *Bi-lift* measure method; finding that you want to evaluate the relationship of $(A \rightarrow B)$ by $\text{lift}(A \rightarrow B)$, you should also study on the

relationship of $\overline{A} \rightarrow B$, so we introduce $\text{lift}(\overline{A} \rightarrow B)$ to adjust $\text{lift}(A \rightarrow B)$. The higher the $\text{lift}(A \rightarrow B)$ is, the better the rule $A \rightarrow B$ is, while the higher the $\text{lift}(\overline{A} \rightarrow B)$ is, the worse the rule $A \rightarrow B$ is. So we propose a *Bi-lift* measure method, $\text{lift}(\overline{A} \rightarrow B)$ as denominator, and $\text{lift}(A \rightarrow B)$ as numerator, namely, ratio of $\text{lift}(A \rightarrow B)$ to $\text{lift}(\overline{A} \rightarrow B)$; *Bi-lift* formula is as follows:

$$\begin{aligned} \text{Bi-lift}(A \rightarrow B) &= \frac{\text{lift}(A \rightarrow B)}{\text{lift}(\overline{A} \rightarrow B)} \\ &= \frac{P(AB)/P(A)P(B)}{P(\overline{A}B)/P(\overline{A})P(B)} \\ &= \frac{P(AB)P(\overline{A})}{P(\overline{A}B)P(A)}. \end{aligned} \quad (11)$$

The premise is $P(\overline{A}B) \neq 0$, and “A” and “B” are not certain event or impossible event. Its value range is $[0, \infty)$. The *Bi-lift* method takes deduction of negative premise as a constraint, to form a bideduction comparing algorithm so as to improve the reliability of the mutual influence between premise and follow-up.

3.2. *Bi-Improve*. Because of the defects of *improve*, the paper put forward *Bi-improve*. Because the probability of former pieces’ occurrence will seriously affect *Improve* evaluation in such a way that when it is high, the *improve* value will be very small all the time. In order to eliminate the influence, we make correction by multiplying the ratio of the occurrence possibility of antecedent to the no occurrence probability of antecedent. *Bi-improve* formula is as follows:

$$\begin{aligned} \text{Bi-improve}(A \rightarrow B) &= [P(B|A) - P(B)] * \frac{P(A)}{P(\overline{A})} \\ &= \frac{P(AB) - P(A)P(B)}{P(\overline{A})}. \end{aligned} \quad (12)$$

Taking Tables 3 and 4, for examples, $\text{Improve}(A' \rightarrow B') = 0.03$ and $\text{Improve}(C \rightarrow D) = 0.11$ only thinking about the *Improve*; rule $(C \rightarrow D)$ is more valuable than rule $(A' \rightarrow B')$. But the fact is also very clear that it can increase “the occurrence possibility of B ” as “A” occurs by up to 39% compared to that when “A” does not occur, while it can increase “the occurrence possibility of event D ” as “C” occurs by up to 19% compared to that when “C” does not occur. So rule $(A' \rightarrow B')$ should be more meaningful than rule $(C \rightarrow D)$. Calculate *Bi-Improve* value through the following formula:

$$\begin{aligned} \text{Bi-improve}(A \rightarrow B) &= [P(B|A) - P(B)] * \frac{P(A)}{P(\overline{A})} \\ &= 0.27, \end{aligned}$$

$$\begin{aligned} & \text{Bi-improve}(C \rightarrow D) \\ &= [P(D | C) - P(D)] * \frac{P(C)}{P(\overline{C})} = 0.225. \end{aligned} \quad (13)$$

$\text{Bi-improve}(A \rightarrow B)$ is higher than $\text{Bi-improve}(C \rightarrow D)$, which accords with the real condition.

3.3. Bi-Confidence. *Confidence* indicates that the appearance of some itemsets will lead to appearance of other itemsets. But we see that the *confidence* of association rules only thinks about the occurrence possibility of “B” when “A” occurs, but not consider the relationship between “A” and “B” when “A” does not occur. So it makes a lot of association rules mining invalid. For the above problems of association rules, we found that the description of *confidence* is not perfect and not enough to show the degree of correlation between itemsets. We put forward the concept of *Bi-confidence*, and its definition is as follows:

$$\begin{aligned} \text{Bi-confidence}(A \rightarrow B) &= \frac{P(AB)}{P(A)} - \frac{P(\overline{A}B)}{P(\overline{A})} \\ &= \frac{P(AB) - P(A)P(B)}{P(A) \times [1 - P(A)]}. \end{aligned} \quad (14)$$

The value range of *Bi-confidence* is $[-1, 1]$. If the *Bi-confidence* value is greater than 0, then $P(AB) > P(A)P(B)$, which shows that “A” and “B” have the positive correlation. If the *Bi-confidence* is equal to 1, then $P(AB) = P(A) = P(B)$, and it shows that “A” and “B” in record set appear together or not. If the *Bi-confidence* is equal to 0, and $P(AB) = P(A)P(B)$, which shows that “A” has no relation with “B.” If the *Bi-confidence* is less than 0, then $P(AB) < P(A)P(B)$; it shows that “A” and “B” have the negative correlation, and negative rules also have research value. *Bi-confidence’s* definition not only contains the correlation factors, but also contains $P(B)$ factor. Therefore, *Bi-confidence* can fully embody the effectiveness of the rules. If we use *Support-Bi-confidence* framework to replace *Support-confidence* framework, it not only can mine association rules effectively, but also can reduce the occurrence of the weak correlation rules.

4. Subjective Interestingness Measures of Association Rules

Generally, the purpose of association rule mining is to obtain certain utility or benefit through the use of some appropriate association rules. So taking the user’s subjective preference or specific application object into consideration, profit targets (or revenue function) are the real key for the users. Thus, two critical problems emerge: firstly, it can bring what kind of utility to the users by the adopted association rules; secondly, the cost or the cost of labor must be taken into consideration when performing the association rules.

4.1. Utility Function: Incremental Monetary Value. The purpose of association rule mining is to translate it into real value and to obtain certain utility. Generally utility function

is the concern point of users. Specifically, utility of rules can be incremental monetary value generated by association rule [20]. Incremental monetary value (*IMV*) is the expected profit (*EP*) under the guidance of rules minus the profit you would expect to receive without the guidance of rules or due to the natural course. Incremental monetary value is defined as follows in that study:

$$\begin{aligned} \text{IMV}(A \rightarrow B) &= [P(B | A) - P(B)] \\ &\quad \times \sum \text{Price}(B_i), \end{aligned} \quad (15)$$

where $\text{Price}(B_i)$ is the unit price of goods B_i and $\sum \text{Price}(B_i)$ is the sum price of all the “B” sets.

Incremental monetary value (*IMV*) has some defects, the probability of former pieces’ occurrence will seriously affect *Improve* evaluation ($\text{Imp.}(A \rightarrow B) = P(B | A) - P(B)$), and, in order to eliminate the influence, we make correction by multiplying by the ratio of the occurrence possibility of antecedent to the no occurrence probability of antecedent. In addition, utility or value should be price minus cost. So we put forward antecedent incremental monetary value (*AIMV*):

$$\begin{aligned} \text{AIMV}(A \rightarrow B) &= \left\{ [P(B | A) - P(B)] \times \frac{P(A)}{P(\overline{A})} \right\} \\ &\quad \cdot \sum [\text{Price}(B_i) - \text{Cost}(B_i)], \end{aligned} \quad (16)$$

where $\text{Price}(B_i)$ is the unit price of goods B_i and $\text{Cost}(B_i)$ is the cost per unit of goods B_i .

4.2. Cost Function. The cost or the cost of labor should be concerned when executing the association rules. Let $\text{Cost}(A \rightarrow B)$ represent the executing cost of rules of ($A \rightarrow B$), such as the unit cost of handling commodity “A” from place A to place B. But the calculation of the so-called executing cost is complicated sometimes. In order to make the study under the maneuverability, the executing cost of rules of ($A \rightarrow B$) can be divided into several levels.

4.3. Interestingness Function Based on Profit (IFBP). On the basis of utility function and the executing cost of rules $\text{Cost}(A \rightarrow B)$, we propose interestingness function based on profit (*IFBP*) with subjective preferences and specific application object. Formula is as follows:

$$\begin{aligned} \text{IFBP}(A \rightarrow B) &= \left\{ [P(B | A) - P(B)] \times \frac{P(A)}{P(\overline{A})} \right\} \\ &\quad \cdot \sum [\text{Price}(B_i) - \text{Cost}(B_i)] \\ &\quad - \text{Cost}(A \rightarrow B). \end{aligned} \quad (17)$$

5. The Measure Framework of Association Rules and Experimental Analysis

5.1. Example Analysis of Objective Interestingness Measures. According to a set of business data in Table 1, take verification

of kinds of measure methods and design measure framework. Because item E appears in all affairs, take it as a kind of certain event and do not think of association rules about E . Set minimum *support* for 20% and minimum *confidence* for 50%. Frequent 2 itemsets calculated through the a priori algorithm are shown in Table 5.

From Table 5, we can see that *support* and *confidence* mainly have the function of basic *support* as classic association rules measure methods, but they can neither distinguish positive and negative correlation nor find the value of various rules. The differences of different measure methods exist, and one single objective interestingness measure method cannot decide which rules are really valuable. That means one must depend on the measure framework which combines multiple indexes.

- (1) The *validity* (Val) is not effective and has huge difference from other measure methods. Each rule in the Measure table of *validity* is valid, but, in fact, many rules have no significance; for example, $F \rightarrow J$, $F \rightarrow R$ are uncorrelated; $F \rightarrow G$, $G \rightarrow N$ have negative inhibitory effect.
- (2) The traditional *support-confidence* framework can eliminate most uncorrelated association rules; but, as low constraints, it can also produce a lot of uncorrelated frequent itemsets and even some negative correlated rules or rules of error.
- (3) *Lift* has a good evaluation result. But obviously *lift* takes events A and B in equivalent position; $lift(A \rightarrow B)$ and $lift(B \rightarrow A)$ are the same; that is to say, if we accept rule $A \rightarrow B$, $B \rightarrow A$ should be also accepted; under such situation, we propose *Bi-lift* to solve it.
- (4) But *Bi-lift* also has a small defect. Its premise is $P(\overline{AB}) \neq 0$, and A and B are not certain event or impossible event. Its value range is $[0, \infty)$.
- (5) *Conviction's* value range is $[0, \infty)$. When the value of *conviction* is "1," it means that "A" does not have relation with "B." And the greater the *conviction* is, the higher the interest of the rule will be. But the *conviction* constraints are too strict, and lots of valuable association rules will be removed. In Table 1, for example, $Conviction(M \rightarrow L) = 0.4 \times 0.8 / 0.5 = 0.64 < 1$; its *conviction* is low, but, in fact, the high interest of M and L may exist. The result of Certainty factor (CF) is very similar to that of *Conviction*.
- (6) Shortcomings of *Improve* ($Imp.$) are obvious. Firstly, how much improvement of probability can be called improvement? Secondly, the probability of former pieces' occurrence will seriously affect *Improve* evaluation in such a way that when it is high, the *improve* value will be very small all the time. Thus, it is difficult to distinguish valuable rules and even tends to make inaccurate value evaluation. Therefore, this paper puts forward *Bi-Improve* aimed at evaluating rules value more accurately.

- (7) The evaluating results of the author's "*Bi-improve*," which is adjusted by the occurrence and nonoccurrence probability of the antecedent, will increase distinction and accuracy.
- (8) *Chi-square analysis* ($Csa.$) is proposed based on the *Improve* ($Imp.$). But the evaluations results indicate that the effects still have not been able to solve the problem of *improve* finally, and its evaluation performance is worse than that of *Bi-improve*.

According to the evaluation results and the performance analysis for measure method, the *validity* ($Val.$) is not effective and sometimes it even appears as essential mistake. For instance, it is sometimes counterproductive, while its evaluation is a promotion.

Though *Improve* ($Imp.$) and *Chi-square analysis* ($Csa.$) do not have the essential mistakes, the stability of their evaluation is not good. Sometimes, they maybe are prone to computation error. However, evaluation results of *Bi-lift*, *Bi-improve*, and *Bi-confidence* are almost the same, and their stabilities of evaluations are high. In conclusion, we can eliminate the *validity* ($Val.$), *improve* ($Imp.$), *Chi-square analysis* ($Csa.$), and other measure indicators. Among the *Support*, *Confidence*, *Lift*, *Bi-lift*, *Bi-improve*, and *Bi-confidence*, seek and build a reasonable measure framework. Procedures are as follow: firstly, use *Support* and *Confidence* threshold to filter out frequent set; secondly, calculate *Bi-lift*, *Bi-improve*, and *Bi-confidence* value; then, according to the *Bi-lift*, *Bi-improve* and the *Bi-confidence* value evaluate association rules comprehensively. Actually the final evaluation results of these three kinds of measure methods are very close and they are perfect, which are shown in Tables 6 and 7. Evaluation results and comparisons of *Bi-lift*, *Bi-improve*, and *Bi-confidence* are shown as Figure 1.

5.2. *Example Analysis Subjective Interestingness Measures.* On the basis of utility function and the executing cost of rules $Cost(A \rightarrow B)$, we propose interestingness function based on profit ($IFBP$) considering subjective preferences and characteristics of specific application object. It plays an important role in the association rules evaluation and selection. According to a set of business data in Table 1, take verification of kinds of measure methods and design measure framework. Because item E appears in all affairs, take it as a kind of certain event and do not think of association rules about E . Set minimum *support* for 30% and minimum *confidence* for 35%. Frequent 2 itemsets calculated through the a priori algorithm are shown in Table 8.

Evaluation results of these five kinds of measure methods are shown in Figure 2. According to the result of evaluation and measure performance analysis, the evaluation results of *Improve* ($Imp.$) and *Chi-square analysis* ($Csa.$) are almost the same, their curves are almost coincidence. And their defects exist obviously. *Lift* and *Bi-improve* have good performance of evaluation and similar tendency. Yet *lift* takes events A and B in equivalence position; if we accept rule $A \rightarrow B$, $B \rightarrow A$ should be also accepted. For this situation, *Bi-improve* is the best choice. On the basis of utility function and the executing cost of rules $Cost(A \rightarrow B)$, we propose

TABLE 5: Various measure evaluation results of various rules.

(a)

Rules	<i>Sup.</i>	<i>Con.</i>	<i>Lift</i>	<i>Bi-lift</i>	<i>Val.</i>	<i>Conv.</i>	<i>Imp.</i>	<i>Bi-imp.</i>	<i>Csa.</i>	<i>CF</i>	<i>Bi-con.</i>
$M \rightarrow J$	0.3	0.75	1.5	2.25	0.1	2	0.25	0.16	1.58	0.50	0.42
$M \rightarrow G$	0.3	0.75	1.08	1.13	0.1	1.2	0.05	0.03	0.35	0.17	0.08
$J \rightarrow G$	0.4	0.8	1.14	1.33	0.1	0.75	0.1	0.1	0.69	0.33	0.2
$I \rightarrow H$	0.3	0.75	1.88	4.5	0.2	2.4	0.35	0.23	2.26	0.58	0.58
$I \rightarrow F$	0.4	1	1.25	1.5	0	/	0.2	0.13	1.58	1.00	0.33
$H \rightarrow F$	0.3	0.75	0.95	0.9	-0.2	0.8	-0.05	-0.03	-0.4	-0.06	-0.08
$R \rightarrow G$	0.5	1	1.42	2.5	0.3	/	0.3	0.3	2.1	1.00	0.6
$N \rightarrow G$	0.4	0.67	0.95	0.89	0.1	0.9	-0.03	-0.05	-0.21	-0.04	-0.08
$R \rightarrow F$	0.4	0.8	1	1	0	1	0	0	0	0.00	0
$N \rightarrow F$	0.5	0.83	1.04	1.11	0.2	1.2	0.03	0.05	0.24	0.15	0.08
$M \rightarrow F$	0.4	1	1.25	1.5	0	/	0.2	0.13	1.58	1.00	0.33
$J \rightarrow F$	0.4	0.8	1	1	0	1	0	0	0	0.00	0
$G \rightarrow F$	0.5	0.71	0.89	0.71	0.2	0.7	-0.09	-0.21	-0.71	-0.11	-0.28
$J \rightarrow M$	0.3	0.6	1.5	3	0.2	1.5	0.2	0.2	1.29	0.33	0.4
$G \rightarrow J$	0.4	0.57	1.14	1.71	0.3	1.17	0.07	0.18	0.44	0.14	0.24
$H \rightarrow I$	0.3	0.75	1.88	4.5	0.2	2.4	0.35	0.23	2.26	0.58	0.58
$F \rightarrow I$	0.4	0.5	1.25	/	0.4	1.2	0.1	0.4	0.65	0.17	0.5
$G \rightarrow R$	0.5	0.71	1.42	/	0.5	1.75	0.21	0.49	1.33	0.42	0.71
$G \rightarrow N$	0.4	0.57	0.95	0.86	0.2	0.93	-0.03	-0.07	-0.19	-0.05	-0.1
$F \rightarrow R$	0.4	0.5	1	1	0.3	1	0	0	0	0.00	0
$F \rightarrow N$	0.5	0.63	1.04	1.25	0.4	1.07	0.03	0.12	0.19	0.08	0.13
$F \rightarrow M$	0.4	0.5	1.25	/	0.4	1.2	0.1	0.4	0.65	0.17	0.5
$F \rightarrow J$	0.4	0.5	1	1	0.3	1	0	0	0	0.00	0
$F \rightarrow G$	0.5	0.63	0.89	0.63	0.3	0.8	-0.07	-0.28	-0.48	-0.10	-0.38
$M \rightarrow L$	0.2	0.5	2.5	/	0.2	0.64	0.3	0.2	2.37	0.38	0.5

(b)

Rules	<i>Lift</i> rank	<i>Bi-lift</i> rank	<i>Val.</i> rank	<i>Conv.</i> rank	<i>Imp.</i> rank	<i>Bi-imp.</i> rank	<i>Csa.</i> rank	<i>CF</i> rank	<i>Bi-con.</i> rank
$M \rightarrow J$	5	9	17	6	5	10	5	6	8
$M \rightarrow G$	14	15	19	12	14	16	14	11	15
$J \rightarrow G$	13	13	18	23	12	14	10	9	13
$I \rightarrow H$	3	6	12	5	2	6	3	5	4
$I \rightarrow F$	11	12	22	3	8	12	7	3	11
$H \rightarrow F$	23	21	25	22	23	21	23	23	22
$R \rightarrow G$	7	8	5	1	3	4	4	2	2
$N \rightarrow G$	22	22	20	20	22	22	22	21	21
$R \rightarrow F$	20	20	24	18	20	20	20	20	20
$N \rightarrow F$	16	16	14	11	15	15	15	14	16
$M \rightarrow F$	10	11	21	2	7	11	6	1	10
$J \rightarrow F$	19	19	23	17	19	19	19	19	19
$G \rightarrow F$	25	24	16	24	25	24	25	25	24
$J \rightarrow M$	4	7	13	8	9	8	9	10	9
$G \rightarrow J$	12	10	6	13	13	9	13	15	12
$H \rightarrow I$	2	5	11	4	1	5	2	4	3
$F \rightarrow I$	9	3	3	10	11	3	12	13	7
$G \rightarrow R$	6	1	1	7	6	1	8	7	1
$G \rightarrow N$	21	23	15	19	21	23	21	22	23
$F \rightarrow R$	18	18	8	16	18	18	18	18	18
$F \rightarrow N$	15	14	4	14	16	13	16	16	14
$F \rightarrow M$	8	2	2	9	10	2	11	12	6
$F \rightarrow J$	17	17	7	15	17	17	17	17	17
$F \rightarrow G$	24	25	9	21	24	25	24	24	25
$M \rightarrow L$	1	4	10	25	4	7	1	8	5

Note: keep two digits after the decimal point.

TABLE 6: The sort of positive association rules.

Rules	<i>Bi-lift</i>	<i>Bi-imp.</i>	<i>Bi-con.</i>	Rank
$G \rightarrow R$	7	0.49	0.71	1
$F \rightarrow M$	6	0.4	0.5	2
$F \rightarrow I$	6	0.4	0.5	3
$R \rightarrow G$	2.5	0.3	0.6	4
$H \rightarrow I$	4.5	0.23	0.58	5
$I \rightarrow H$	4.5	0.23	0.58	6
$M \rightarrow L$	5	0.2	0.5	7
$J \rightarrow M$	3	0.2	0.4	8
$G \rightarrow J$	1.71	0.18	0.24	9
$M \rightarrow J$	2.25	0.16	0.42	10
$M \rightarrow F$	1.5	0.13	0.33	11
$I \rightarrow F$	1.5	0.13	0.33	12
$F \rightarrow N$	1.25	0.12	0.13	13
$J \rightarrow G$	1.33	0.1	0.2	14
$N \rightarrow F$	1.11	0.045	0.08	15
$M \rightarrow G$	1.13	0.03	0.08	16

TABLE 7: The sort of negative association rules and meaningless rules.

Rules	<i>Bi-lift</i>	<i>Bi-imp.</i>	<i>Bi-con.</i>	Rank
$H \rightarrow F$	0.9	-0.03	-0.08	17
$N \rightarrow G$	0.89	-0.045	-0.08	18
$G \rightarrow N$	0.86	-0.07	-0.1	19
$G \rightarrow F$	0.71	-0.21	-0.28	20
$F \rightarrow G$	0.63	-0.28	-0.38	21
$F \rightarrow J$	1	0	0	22
$F \rightarrow R$	1	0	0	23
$J \rightarrow F$	1	0	0	24
$R \rightarrow F$	1	0	0	25

interestingness function based on profit (*IFBP*) considering subjective preferences and characteristics specific application object. It plays an important role in the association rules evaluation and selection.

6. Conclusion

As a statistics based method, association rule mining has certain limitations. First of all, the generation of the association rules is totally based on the fact data without considering the relationship between the rules. Secondly, affected by data quality and selection of threshold, the generator may produce useless rules or even lose some useful rules. Thirdly, the expression ability of association rules is limited. Thus, evaluating the reliability of the obtained association rules becomes one of the hot spots for researchers. Study on traditional association rules mining is based on *support-confidence* framework, and the rules are called strong association rules only when they satisfy both thresholds of *support* and *confidence*. However, sometimes strong association rules are not what users are interested in and are even misleading.

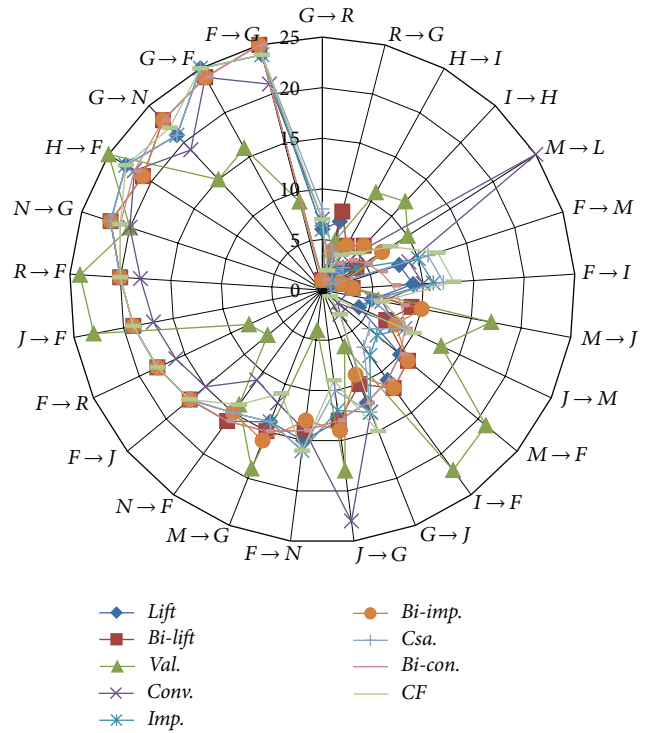


FIGURE 1: The radar map of evaluation results of measure methods.

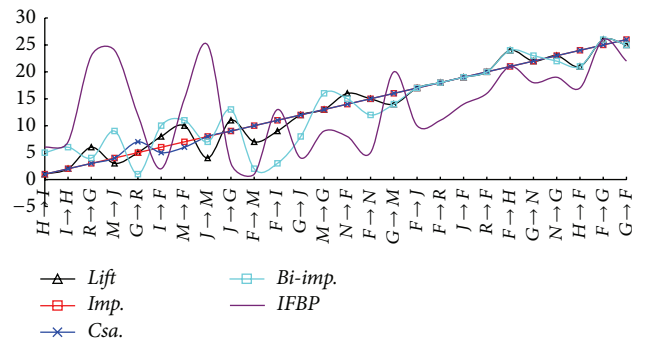


FIGURE 2: Evaluation results of these five kinds of measure methods (*Lift*, *Imp.*, *Csa.*, and *N-Imp.*, *IFBP*).

Thus, to further analyze and evaluate the mined rules in order to find the most valuable association.

Generally, since most rules with high *support* are obvious or are already known by users, low *support* rules that provide users with some interesting new knowledge may be more novel than high *support* rules. However, with too low *support* threshold, it can also produce the combination explosion problem. So the best way to resolve this dilemma is to first set a low *support* threshold or use dynamic *support* threshold to complete a series of mining and then employ the new association rules measure framework to screen mining results and extract the most valuable and interesting association rules at the same time.

TABLE 8: Various measure evaluation results of various rules.

(a)

Rules	Cost	Profit	Sup.	Con.	Lift	Imp.	Csa.	Bi-imp.	IFBP
$H \rightarrow I$	1	2	0.3	0.75	1.88	0.35	2.26	0.23	0.23
$I \rightarrow H$	5	6	0.3	0.75	1.88	0.35	2.26	0.23	0.23
$R \rightarrow G$	5	4	0.5	1	1.42	0.3	2.1	0.3	-0.3
$M \rightarrow J$	7	5	0.3	0.75	1.5	0.25	1.58	0.16	-0.32
$G \rightarrow R$	3	3	0.5	0.71	1.42	0.21	1.33	0.49	0
$I \rightarrow F$	1	5	0.4	1	1.25	0.2	1.58	0.13	0.52
$M \rightarrow F$	5	5	0.4	1	1.25	0.2	1.58	0.13	0
$J \rightarrow M$	6	3	0.3	0.6	1.5	0.2	1.29	0.2	-0.6
$J \rightarrow G$	2	7	0.4	0.8	1.14	0.1	0.69	0.1	0.5
$F \rightarrow M$	1	3	0.4	0.5	1.25	0.1	0.65	0.4	0.8
$F \rightarrow I$	2	2	0.4	0.5	1.25	0.1	0.65	0.4	0
$G \rightarrow J$	3	5	0.4	0.57	1.14	0.07	0.44	0.18	0.36
$M \rightarrow G$	4	7	0.3	0.75	1.08	0.05	0.35	0.03	0.09
$N \rightarrow F$	1	5	0.5	0.83	1.04	0.03	0.24	0.045	0.18
$F \rightarrow N$	2	4	0.5	0.63	1.04	0.03	0.19	0.12	0.24
$G \rightarrow M$	5	3	0.3	0.43	1.08	0.03	0.19	0.07	-0.14
$F \rightarrow J$	3	5	0.4	0.5	1	0	0	0	0
$F \rightarrow R$	1	3	0.4	0.5	1	0	0	0	0
$J \rightarrow F$	2	5	0.4	0.8	1	0	0	0	0
$R \rightarrow F$	2	5	0.4	0.8	1	0	0	0	0
$F \rightarrow H$	5	7	0.3	0.38	0.95	-0.02	-0.13	-0.08	-0.16
$G \rightarrow N$	3	4	0.4	0.57	0.95	-0.03	-0.19	-0.07	-0.07
$N \rightarrow G$	2	4	0.4	0.67	0.95	-0.03	-0.21	-0.045	-0.09
$H \rightarrow F$	3	5	0.3	0.75	0.95	-0.05	-0.4	-0.03	-0.06
$F \rightarrow G$	1	4	0.5	0.63	0.89	-0.07	-0.48	-0.28	-0.84
$G \rightarrow F$	4	5	0.5	0.71	0.89	-0.09	-0.71	-0.21	-0.21

(b)

Rules	Cost	Profit	Lift Rank	Imp. Rank	Csa. Rank	Bi-imp. Rank	IFBP Rank
$H \rightarrow I$	1	2	1	1	1	5	6
$I \rightarrow H$	5	6	2	2	2	6	7
$R \rightarrow G$	5	4	6	3	3	4	23
$M \rightarrow J$	7	5	3	4	4	9	24
$G \rightarrow R$	3	3	5	5	7	1	12
$I \rightarrow F$	1	5	8	6	5	10	2
$M \rightarrow F$	5	5	10	7	6	11	15
$J \rightarrow M$	6	3	4	8	8	7	25
$J \rightarrow G$	2	7	11	9	9	13	3
$F \rightarrow M$	1	3	7	10	10	2	1
$F \rightarrow I$	2	2	9	11	11	3	13
$G \rightarrow J$	3	5	12	12	12	8	4
$M \rightarrow G$	4	7	13	13	13	16	9
$N \rightarrow F$	1	5	16	14	14	15	8
$F \rightarrow N$	2	4	15	15	15	12	5
$G \rightarrow M$	5	3	14	16	16	14	20
$F \rightarrow J$	3	5	17	17	17	17	10
$F \rightarrow R$	1	3	18	18	18	18	11
$J \rightarrow F$	2	5	19	19	19	19	14
$R \rightarrow F$	2	5	20	20	20	20	16
$F \rightarrow H$	5	7	24	21	21	24	21
$G \rightarrow N$	3	4	22	22	22	23	18
$N \rightarrow G$	2	4	23	23	23	22	19
$H \rightarrow F$	3	5	21	24	24	21	17
$F \rightarrow G$	1	4	26	25	25	26	26
$G \rightarrow F$	4	5	25	26	26	25	22

Note: keep two digits after the decimal point.

Conflict of Interests

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

Acknowledgments

Research is supported by The National Key Technology R and D Program of China (Grant 2014BAH24F06), Natural Science Foundation of Zhejiang Province (Grants LY14F020002, LY15G010001), and the Contemporary Business and Trade Research Center of Zhejiang Gongshang University which is the Key Research Institutes of Social Sciences and Humanities Ministry of Education.

References

- [1] P. N. Tan, M. Steinbach, and V. Kumar, *Introduction to Data Mining*, Pearson—Addison-Wesley, Boston, Mass, USA, 1st edition, 2005.
- [2] J. H. Fowler and M. Laver, “A tournament of party decision rules,” *Journal of Conflict Resolution*, vol. 52, no. 1, pp. 68–92, 2008.
- [3] R. Agrawal and R. Srikant, “Fast algorithms for mining association rules,” in *Proceedings of the 20th International Conference on Very Large Data Bases (VLDB '94)*, pp. 487–499, 1994.
- [4] R. Srikant and R. Agrawal, “Mining generalized association rules,” in *Proceedings of the 21st International Conference on Very Large Data Bases*, pp. 407–418, Zurich, Switzerland, 1995.
- [5] M.-C. Tseng, W.-Y. Lin, and R. Jeng, “Incremental maintenance of generalized association rules under taxonomy evolution,” *Journal of Information Science*, vol. 34, no. 2, pp. 174–195, 2008.
- [6] M.-C. Tseng and W.-Y. Lin, “Maintenance of generalized association rules with multiple minimum supports,” *Intelligent Data Analysis*, vol. 8, no. 4, pp. 417–436, 2004.
- [7] M.-C. Chen, “Ranking discovered rules from data mining with multiple criteria by data envelopment analysis,” *Expert Systems with Applications*, vol. 33, no. 4, pp. 1110–1116, 2007.
- [8] M. Toloo, B. Sohrabi, and S. Nalchigar, “A new method for ranking discovered rules from data mining by DEA,” *Expert Systems with Applications*, vol. 36, no. 4, pp. 8503–8508, 2009.
- [9] L. Geng and H. J. Hamilton, “Interestingness measures for data mining: a survey,” *ACM Computing Surveys*, vol. 38, no. 3, pp. 1–32, 2006.
- [10] N. Hoque, B. Nath, and D. K. Bhattacharyya, “A new approach on rare association rule mining,” *International Journal of Computer Applications*, vol. 53, no. 3, pp. 1–6, 2012.
- [11] J. Zhang, Y. Wang, and J. Feng, “Attribute index and uniform design based multiobjective association rule mining with evolutionary algorithm,” *The Scientific World Journal*, vol. 2013, Article ID 259347, 16 pages, 2013.
- [12] J. Han and M. Kamber, *Data Mining Concept and Technology*, China Machine Press, Beijing, China, 2007.
- [13] G. Piatetsky-Shapiro, “Discovery, analysis, and presentation of strong rules,” in *Knowledge Discovery in Databases*, G. Piatetsky-Shapiro and W. Frawley, Eds., pp. 229–238, AAAI/MIT Press, 1991.
- [14] L. Ma and W. Jie, “Research on judgment criterion of association rules,” *Control and Decision*, vol. 18, no. 3, pp. 277–280, 2003.
- [15] S. Brin, R. Motwani, J. D. Ullman, and S. Tsur, “Dynamic itemset counting and implication rules for market basket data,” in *Proceedings of the ACM SIGMOD International Conference on Management of Data (SIGMOD '97)*, pp. 255–264, ACM, Tucson, Ariz, USA, May 1997.
- [16] P.-N. Tan, V. Kumar, and J. Srivastava, “Selecting the right objective measure for association analysis,” *Information Systems*, vol. 29, no. 4, pp. 293–313, 2004.
- [17] Y. Li, C. Wu, and K. Wang, “A new interestingness measures for Ming association rules,” *Journal of the China Society for Scientific and Technical Information*, vol. 30, no. 5, pp. 503–507, 2011.
- [18] J. Chen and Y. Gao, “Evaluation criterion for association rules with influence degree,” *Computer Engineering and Application*, vol. 45, no. 8, pp. 141–143, 2009.
- [19] F. Berzal, I. Blanco, D. Sanchez, and M.-A. Vila, “Measuring the accuracy and interest of association rules: a new framework,” *Intelligent Data Analysis*, vol. 6, no. 3, pp. 221–235, 2002.
- [20] D. H. Choi, B. S. Ahn, and S. H. Kim, “Prioritization of association rules in data mining: multiple criteria decision approach,” *Expert Systems with Applications*, vol. 29, no. 4, pp. 867–878, 2005.



Hindawi

Submit your manuscripts at
<http://www.hindawi.com>

