# Toward Multidatabase Mining: Identifying Relevant Databases

Huan Liu, *Senior Member*, *IEEE*, Hongjun Lu, and Jun Yao

**Abstract**—Various tools and systems for knowledge discovery and data mining are developed and available for applications. However, when we are immersed in heaps of databases, an immediate question is where we should start mining. It is not true that the more databases, the better for data mining. It is only true when the databases involved are relevant to a task at hand. In this paper, breaking away from the conventional data mining assumption that many databases be joined into one, we argue that the first step for multidatabase mining is to identify databases that are most likely relevant to an application; without doing so, the mining process can be lengthy, aimless, and ineffective. A measure of relevance is thus proposed for mining tasks with an objective of finding patterns or regularities about certain attributes. An efficient algorithm for identifying relevant databases is described. Experiments are conducted to verify the measure's performance and to exemplify its application.

**Index Terms**—Multiple databases, data mining, query, relevance measure.

✦

## 1 INTRODUCTION

WITH more and more databases created, an increasingly pressing issue is how to make efficient use of them. To address the issue, there has been a recent surge of research interest on knowledge discovery and data mining [1], [10], [7], [33]. While researchers are trying to develop efficient algorithms to cope with large volumes of data, little work has been devoted to the *data aspect* in the knowledge discovery process. In most organizations, data is rarely specially collected and stored for the purpose of mining knowledge, but usually as the byproducts of other tasks [32]. Furthermore, with the development of technologies, it is not uncommon that an organization has a large number of database systems and diverse data sources.

Although most data mining algorithms assume a single data set, for real world applications, practitioners have to face the problem of discovering knowledge from multiple databases. In order to do so, one way is to employ a brute force approach to join the available tables into a single large table upon which existing data mining techniques or tools can be applied. There are several problems for this approach in real world applications. First, database integration itself is still a problematic area, especially where the source domains differ. Second, all tables with foreign key references need to be joined together to produce a single combined table. The size of the resulting table, in terms of both the number of records and the number of attributes, will be much larger than the original individual tables. The

increase of data size not only prolongs the running time of mining algorithms, but also affects the behavior of mining algorithms. From the viewpoint of statistics, joining one relevant database with an irrelevant one will result in a more difficult task to find useful patterns as search space is enlarged by irrelevant attributes. For example, there are two binary-valued databases and each has $N$ attributes, assuming that one database is irrelevant and $N/2$ attributes can be found in both databases. Simply working on the relevant database, the hypothesis space is $2^{2^N}$; after joining, it is $2^{2^{3N/2}}$. For this simplified analysis, we have not yet considered the factor of missing values due to joining. This factor will certainly increase the difficulty of data mining, too. Third, if databases are joined and data mining algorithms are applied, the users face the problem of identifying interesting patterns from a large number of discovered rules. In practice, it is too easy to discover a huge number of patterns in a database [20], [26], [28]; however, it is difficult for users to search in all the discovered patterns for useful ones. The redundant, useless, or uninteresting patterns can be even more easily generated when there are quite a number of databases irrelevant to the mining task. Therefore, as in any effective knowledge discovery process, the first important step in mining multiple databases is indeed to select those databases that are *relevant* to a specific mining task.

In this paper, we will address the *relevance* problem: identifying databases relevant to a particular data mining task in multiple databases. Without loss of generality, we call each database a relation or a table and assume that 1) a specific mining task is related to the property of certain attributes, which can be expressed using query predicates, and 2) the higher order correlations of attributes with a query are at least partially reflected in their first order correlations. It is important to note that, unlike the conventional query processing, our task is to select the databases relevant to a given query predicate rather than identifying the databases matching the query.

- *H. Liu is with the Department of Computer Science and Engineering, Arizona State University, PO Box 875406, Tempe, AZ 85287-5406. E-mail: hliu@asu.edu.*
- *H. Lu is with the Department of Computer Science, Hong Kong University of Science and Technology, Clear Water Bay, Kowloon, Hong Kong, China. E-mail: luhj@cs.ust.hk.*
- *J. Yao is with Mokonet Internet Inc., 40-31 3FL, 68th St., Woodside, NY 11377. E-mail: yaojun@excite.com.*

TABLE 1
Database 1 of Diet Habits

| ID | group | alcohol | on-diet | snack-between-meals | favorite-non-veg |
|----|-------|---------|---------|---------------------|------------------|
| 950351 | Chinese | never | no | sometimes | fish |
| 950301 | Chinese | never | no | seldom | pork |
| 950282 | Chinese | sometimes | no | seldom | fish |
| 940112 | Chinese | often | no | often | chicken |
| 940023 | Chinese | sometimes | no | sometimes | beef |
| 938976 | Chinese | sometimes | no | sometimes | egg |
| 950612 | Russian | never | no | seldom | beef |
| 950122 | Russian | often | no | sometimes | beef |
| 940227 | Russian | sometimes | no | sometimes | chicken |
| 938567 | Russian | sometimes | no | often | pork |
| 950348 | Indian | often | no | sometimes | fish |
| 950312 | Indian | sometimes | no | seldom | pork |
| 950123 | Indian | never | no | sometimes | chicken |
| 940247 | Indian | sometimes | no | sometimes | fish |
| 940100 | Indian | never | no | sometimes | beef |

The problem of identifying relevant databases can be stated as follows:

> Given $n$ data tables (relations), $D_k(1 \leq k \leq n)$, each of which consists of a number of attributes. A query predicate $Q$ is expressed in the form of "$A_i$ relop $C$" where $A_i$ is an attribute, relop is one of the relational operators $(=, >, <, \leq, \geq, \neq)$, and $C$ is a value in its domain. Identifying relevant databases is to select relevant tables (relations) that contain specific, reliable, and statistically significant information pertaining to the query.

By being reliable and significant, we mean that the information is consistent and accurate to a certain degree and not due to sheer chance, even when its occurrence is not frequent.

We start mining tasks with query predicates for two reasons: 1) Although users may not exactly know what kind of knowledge to mine, they usually have certain rough ideas. Queries are easy to form and they are most familiar and natural to database users: The user's interest, prior knowledge, and intention can be conveyed by means of queries. 2) Among various data mining tasks (classification, characterization, generalization, association, etc.), most of them aim at discovering knowledge that can be expressed as relationships among values of attributes [4], [23], [11], [27], [30], [8]. For example, classification rules, characteristic rules, and association rules can be expressed in the form of $A \rightarrow B$, where both $A$ and $B$ are conjunctions of attribute values. Queries can easily take care of these forms of knowledge.

There are two key issues concerning solutions to the selection problem: 1) The solution should be efficient in using databases and 2) only the relevant databases should be chosen for mining tasks. Note that it is not our intention to propose yet another data mining method [31]. What we are concerned about is having a relatively simple and fast method by which we can identify the databases (tables) that are relevant to a mining task. After that, various existing techniques can be used to conduct mining tasks. Major contributions of our work include:

- Illustrating the importance and necessity of identifying relevant databases when mining multiple databases;
- Establishing a quantitative measure that allows us to identify relevant database (tables) before mining; and
- Developing an efficient implementation of computing the relevance measure.

The remainder of the paper is organized as follows: Section 2 gives a motivating example to illustrate the issue of relevance. A relevance measure that can tell the relevant tables from the irrelevant ones is presented in Section 3. Section 4 describes an algorithm that identifies relevant databases using the proposed measure. Section 5 presents some preliminary results of using the measure on a real database. Section 6 discusses the relationship between our work and some related work. Finally, Section 7 concludes the paper.

## 2   A MOTIVATING EXAMPLE

Before we design the relevance measure for identifying relevant databases, let's look at a simple example. We will use two small databases to argue that, for data mining tasks involving multiple databases, it is better to identify a relevant database first before applying mining techniques.

### 2.1   Two Artificial Databases

Our example databases are about people's diet habits. It is the user's intention to discover some *useful* knowledge about Chinese diet habits from the databases. Here, and in most KDD applications, being useful means something that should potentially lead to some useful actions by a user [6]. Specific information is one type of usefulness. In our example, the useful knowledge should be some diet habits specific to Chinese. Two tables (databases) with desired attributes are as shown in Tables 1 and 2. Let us look at Table 1 first. From the attribute values summarized in Table 3, we observe the following:

TABLE 2
Database 2 of Diet Habits

| ID | group | main-food | regular eating times | drink | vegetarian |
|----|-------|-----------|----------------------|-------|------------|
| 950578 | Chinese | rice | 3 | tea | no |
| 950351 | Chinese | rice | 3 | tea | no |
| 950301 | Chinese | rice | 3 | tea | no |
| 950282 | Chinese | rice | 3 | tea | no |
| 940226 | Chinese | rice | 3 | cola | no |
| 940112 | Chinese | rice | 3 | tea | yes |
| 950612 | Russian | bread | 3 | coffee | no |
| 950122 | Russian | bread | 3 | cola | no |
| 940227 | Russian | rice | 3 | cola | yes |
| 940121 | Russian | bread | 3 | tea | no |
| 950348 | Indian | bread | 3 | coffee | no |
| 950312 | Indian | bread | 3 | coffee | yes |
| 950123 | Indian | rice | 3 | cola | no |
| 940247 | Indian | rice | 3 | coffee | no |
| 940109 | Indian | bread | 3 | tea | no |

- No person is on a diet, regardless of his ethnic group;
- For both Chinese and Russian, 50 percent of records show that they take alcohol sometimes;
- Similarly, 50 percent of Chinese and Russian have snacks between regular meals; and
- It seems difficult to arrive at a convincing conclusion on a favorite nonvegetable food since a very small number of tuples are available: There are five different foods listed and the largest group only has six records available.

From the above observations, we can see that the database in fact does not contain information specific to Chinese. In such a case, it is fair to say that the database is irrelevant with respect to the query about Chinese.

TABLE 3
Analysis of Data in Database 1

| group | | Chinese | Russian | Indian | Total |
|-------|--|---------|---------|--------|-------|
| No Tuples | | 6 | 4 | 5 | 15 |
| alcohol | never | 2 | 1 | 1 | 4 |
| | sometimes | 3 | 2 | 2 | 7 |
| | often | 1 | 1 | 2 | 4 |
| on-diet | no | 6 | 4 | 5 | 15 |
| snack-between-meals | seldom | 2 | 1 | 1 | 4 |
| | sometimes | 3 | 2 | 4 | 9 |
| | often | 1 | 1 | 0 | 2 |
| favorite-non-veg. | pork | 1 | 1 | 1 | 3 |
| | chicken | 1 | 1 | 1 | 3 |
| | beef | 1 | 2 | 1 | 4 |
| | fish | 2 | 0 | 2 | 4 |
| | egg | 1 | 0 | 0 | 1 |

It is easy to see that the database in Table 2 is relevant to the query about Chinese since we can at least derive such a statement (rule):

*"If the group is Chinese, the main food is rice"*

since all records about Chinese show that the main food is rice, whereas records about Russian and Indian do not indicate such a habit. We have seen that, although both databases contain information about the diet habits of Chinese, one of them contains relevant information, but the other does not.

## 2.2 Possible Solutions for Mining Multiple Databases

We were able to manually pick one database over the other in the above simple example because the databases are small. We need a method for large databases for selection. Let's examine the existing method and check if some of them can serve the purpose of selection, pretending no knowledge about which database contains interesting information.

1. The first solution is to combine two tables together, i.e., joining them on the common attribute *ID* before applying mining algorithms. There are some problems in doing so. a) The resulting tuples will have 10 attributes besides *ID*, i.e., twice as many as in the original databases. b) The resulting table will have more tuples. In our example, the number of tuples increases to 19 (four more than the original one). And, c) for certain tuples, some attributes will have missing values since they only appear in one of the databases. All these will surely make the subsequent mining task more complex.

2. The second solution could be issuing database queries to retrieve information from the databases and see which one should be focused on. In our example, if an SQL selection query with condition "**WHERE** group = Chinese" is posed, both databases will return six tuples. There is no indication whether one of them is relevant or not.

3. We can apply some data mining algorithms to the two databases; some patterns will be generated from the irrelevant database. For instance, using the concept of *support* and *confidence* level in mining for association rules [9], we can have rules like:

$$group = \text{Chinese} \rightarrow \ on\text{-}diet = \text{no}$$

with 40 percent support and 100 percent confidence. We do not think such information is really interesting since *on-diet* = no in fact holds for all tuples. In other words, what we found is not specific to Chinese. We can also apply certain classification algorithms, such as CART, CN2, and C4.5 [3], [4], [27], to each database, specifying attribute *group* as the class label. Using C4.5, for example, we can induce classification rules from both databases whose accuracy rates are better than that of just choosing the most probable class label. Therefore, based on whether any rules can be induced, we cannot distinguish the relevant database from the irrelevant one.

From the above, commonly used methods do not seem to work for database selection or efficiently distinguishing the irrelevant database from the relevant one. We need a new method to identify relevant databases first and to facilitate efficient mining multiple databases later.

# 3 RELEVANCE MEASUREMENT

In this section, we discuss a quantitative measure of relevance. The relevant databases identified by the measure should contain specific information pertaining to our mining task, which is reliable and statistically significant.

## 3.1 Relevance Factor of Selector

We call "*A relop C*" a *selector*, where $A$ is an attribute name which is not referenced by the query predicate $Q$, $relop \in \{=, <, >, \leq, \geq, \neq\}$, and $C$ is a constant value in the domain of $A$. We define the *relevance factor* of selector $s$, such as "drink = tea," with respect to $Q$, such as "group = Chinese," below:

$$RF(s, Q) = \Pr(s|Q) \Pr(Q) \log \frac{\Pr(s|Q)}{\Pr(s)}, \qquad (1)$$

where $\Pr(Q)$ and $\Pr(s)$ are priors and estimated by the ratios of how frequently they appear in a database and $\Pr(s|Q)$ is the posterior about the frequency ratio of $s$ appearing given that $Q$ occurs. The rationale of defining relevance as in (1) is that $\frac{\Pr(s|Q)}{\Pr(s)}$ shows the degree of the deviation of the posterior from the prior. This ratio tells us different relationships between $\Pr(s|Q)$ and $\Pr(s)$.

**Case 1:** If $\frac{\Pr(s|Q)}{\Pr(s)}$ is close to 1, i.e., $\Pr(s|Q) \simeq \Pr(s)$, $s$ is independent of $Q$.

**Case 2:** If $\frac{\Pr(s|Q)}{\Pr(s)}$ is close to 0, i.e., $\Pr(s|Q)$ is almost 0, $s$ rarely occurs given $Q$.

**Case 3:** If $\frac{\Pr(s|Q)}{\Pr(s)}$ is less than 1, $s$ is not used frequently enough when $Q$ is given.

**Case 4:** If $\frac{\Pr(s|Q)}{\Pr(s)}$ is greater than 1, then $s$ occurs more often given $Q$ than without $Q$, hence, $s$ and $Q$ are correlated.

It is clear that, in the context of finding relevant databases, we are only interested in the last case. Therefore, the logarithm function is taken since $\log \frac{\Pr(s|Q)}{\Pr(s)}$ is either 0 or less than 0 for Cases 1 and 3. Adding a weight $\Pr(s|Q)$ to the measure, we make sure that, in Case 2, $RF$ is also close to 0. $\Pr(s|Q)$ also indicates how valid the correlation between $s$ and $Q$ is. Another weight $\Pr(Q)$ is introduced to take into account how frequently $Q$ occurs. This is because:

**Case 5:** If $\Pr(Q)$ is close to 0, then $Q$ seldom occurs in the database. Hence, it is not statistically significant to conclude that something is related to $Q$.

Thus, we develop a measure $RF$ that is able to distinguish all five cases. The higher the value $RF$, the more relevant a selector $s$ with respect to query $Q$. If $RF$ for all selectors of a database is close to or less than 0, the database is irrelevant to $Q$. The changes in the proportion of irrelevant information can also be seen from the $RF$ definition (1). Assuming that $\Pr(Q)$ and $\Pr(s)$ are fixed, we can see that if $Q$ and $s$ are positively relevant (i.e., $\Pr(s|Q) > \Pr(s)$), $RF$ is positive; if $Q$ and $s$ are negatively relevant (i.e., $\Pr(s|Q) < \Pr(s)$); if $Q$ and $s$ are irrelevant (i.e., independent of each other), $RF$ is 0.

## 3.2 Relevant Databases

With the above definition about the relevance factor of selectors, we can have the following definition about the relevance of databases.

**Definition 1.** *A selector $s$ is relevant to $Q$ if $RF(s, Q) > \delta$, where $\delta (> 0)$ is the given threshold. A table is relevant to $Q$ if there exists at least one selector $s_j$, $A_i$ relop $C$, where $s_j$ is relevant to $Q$.*

Now, we have a quantitative approach to distinguishing irrelevant databases from the relevant ones. Let us compute the values of $RF$ for the example databases in the previous section. Given *Q: Group* = Chinese, Tables 4 and 5 list the $RF$ values for all the selectors in Tables 1 and 2, respectively. $\Pr(s|Q)$ is also given as reference.

From Table 4, we can see that the maximum $RF$ value is 0.088. If we set $\delta$ to 0.1, the database will be considered as irrelevant since there does not exist any selector relevant to $Q$.

Let's check Table 5. With the same value of $\delta$, two relevant selectors, *main food = rice* and *drink = tea*, are considered relevant to $Q$. We conclude that the database is relevant. From the two simple illustrative examples, we can see that the measure $RF$ can tell a relevant database from an irrelevant one.

In order to determine whether a database is relevant to a given query predicate $Q$, we need to examine every attribute and values. For categorical attributes, selectors can be formed for each distinct value. For noncategorical attributes, we can either employ the idea of ranges or discretize attribute values [19], which divides attribute values into a limited number of intervals and maps values in each interval to an integer. When discretization is

TABLE 4
Relevance Factors of *RF* in Database One

| No | Selector $s$ | $RF(s, Q)$ | $Pr(s|Q)$ |
|----|----|----|----|
| 1 | alcohol=never | 0.000000 | 0.333333 |
| 2 | alcohol=sometimes | 0.019907 | 0.500000 |
| 3 | alcohol=often | -0.017536 | 0.166667 |
| 4 | on diet=no | 0.000000 | **1.000000** |
| 5 | snacks between meal=sometimes | -0.052607 | 0.500000 |
| 6 | snacks between meal=seldom | 0.042924 | 0.333333 |
| 7 | snacks between meal=often | 0.021462 | 0.166667 |
| 8 | favorite non-veg=fish | 0.042924 | 0.333333 |
| 9 | favorite non-veg=pork | -0.017536 | 0.166667 |
| 10 | favorite non-veg=chicken | -0.017536 | 0.166667 |
| 11 | favorite non-veg=beef | -0.045205 | 0.166667 |
| 12 | favorite non-veg=egg | **0.088129** | 0.166667 |

needed, we resort to Chi2—a discretization system based on the $\chi^2$ statistics and an inconsistency measure [17]. The system is available for research purposes upon request.

### 3.3 Choosing the Threshold

One issue in using the relevance measure is how to choose the threshold $\delta$. In this section, we study in detail how to choose threshold value for $RF$. The definition of $RF$ consists of two parts: $\Pr(s|Q) \times \Pr(Q)$—the term before the logarithm and $\Pr(s|Q)/\Pr(s)$—the logarithmic term. As mentioned before, $\Pr(s|Q)$ indicates how strong the correlation between $s$ and $Q$ is. By multiplying $\Pr(s|Q)$ by $\Pr(Q)$, the first part is in fact $\Pr(Q \wedge s)$, i.e., the ratio between the number of tuples for which both $Q$ and $s$ hold and the total number of tuples in the database. The second part $\Pr(s|Q)/\Pr(s)$, which can be estimated using the ratio of the number of tuples for which both $s$ and $Q$ hold and the number of tuples where $s$ holds, reflects specificity of $s$ to $Q$ with respect to the whole database. In our definition, the

selectors with larger values for both parts will be considered more specific to a query. In Fig. 1, we draw $\Pr(Q \wedge s) \times \log \Pr(s|Q)/\Pr(s) = \delta_i$ for $\delta_i = 0.05, 0.1, 0.2$, and $0.5$ using $\Pr(Q \wedge s)$ and $\Pr(s|Q)/\Pr(s)$ as variables for the y-axis and x-axis, respectively. For a threshold value $\delta_i$ to satisfy $RF = \Pr(Q \wedge s) \times \log \Pr(s|Q)/\Pr(s) \geq \delta_i$, the values of $\Pr(Q \wedge s)$ and $\Pr(s|Q)/\Pr(s)$ should be in the right-upper part of the curve ($y = \delta_i/\log x$). In other words, a database is identified as relevant only if it contains at least one selector $s$ such that $\Pr(Q \wedge s)$ and $\Pr(s|Q)/\Pr(s)$ fall into the right-upper region of the curve in Fig. 1. Note that the value of $\Pr(s|Q)/Pr(s)$ which reflects the degree of deviation of the posterior from the prior could be very large. However, $\Pr(Q \wedge s)$ is usually not very high in real databases, especially in large databases. The value of $RF$ will decrease rapidly as the value of $\Pr(Q \wedge s)$ decreases. So, we need to lower the threshold when database size is large.

For most data mining tasks, $\Pr(Q \wedge s)$, the percentage of tuples for which both $Q$ and $s$ hold, should not be too low since this value indicates the relative significance in terms of statistics. However, it is still unclear what the lower bound of $\Pr(Q \wedge s)$ is. In order to make sure that $\Pr(Q \wedge s)$ is not too small, the query predicate should cover a sufficient number of tuples in the database in the first place, i.e., $\Pr(Q)$ should not be too small. Usually, we know the total number of tuples in a database. Therefore, it is relatively easy to specify the lower bound of $\Pr(Q)$ (the percentage of tuples covered by a query). We can show that, with the estimated $\Pr(Q)$, we can obtain an upper bound for $RF$.

**Theorem 1.** *Given* $\Pr(Q)$, *the following holds:*

$$RF \leq \Pr(Q) \log(1/\Pr(Q)) \qquad (2)$$

TABLE 5
Relevance Factors of *RF* in Database Two

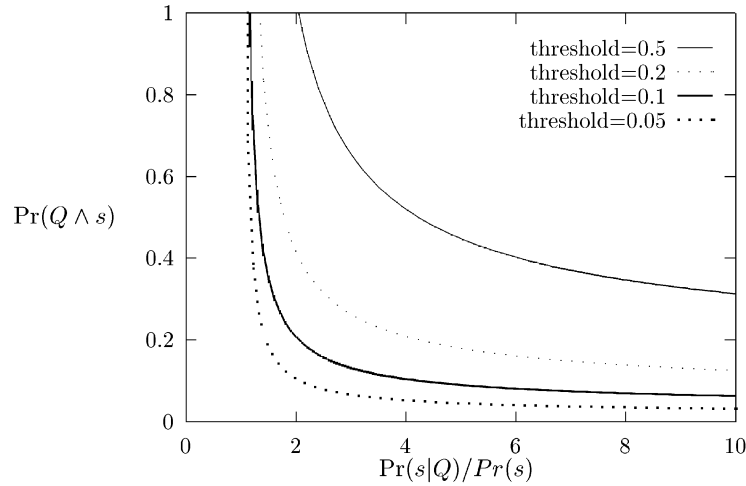| No | Selector $s$ | $RF(s, Q)$ | $Pr(s|Q)$ |
|----|----|----|----|
| 1 | main food=rice | **0.294786** | **1.000000** |
| 2 | regular eating times=3 | 0.000000 | **1.000000** |
| 3 | drink=tea | **0.278834** | 0.833333 |
| 4 | drink=cola | -0.045205 | 0.166667 |
| 5 | vegetarian=no | 0.019631 | 0.833333 |
| 6 | vegetarian=yes | -0.017536 | 0.166667 |

Fig. 1. $\Pr(Q \wedge s)$ as the function of $\Pr(s|Q)/\Pr(s)$ on different threshold values.

**Proof.** Let us define $\Pr(Q) = \frac{N_q}{N}$; $\Pr(s) = \frac{N_s}{N}$; $\Pr(s|Q) = \frac{N_{s \wedge q}}{N_q}$, then

$$RF = \Pr(Q \wedge s) \log \left( \frac{N_{s \wedge q}}{N_q} \times \frac{N}{N_s} \right)$$

$$= \Pr(Q \wedge s) \log \left( \frac{N_{s \wedge q}}{N_s} \times \frac{N}{N_q} \right)$$

$$= \Pr(Q \wedge s) \log \left( \frac{N_{s \wedge q}}{N_s} \times \frac{1}{\Pr(Q)} \right).$$

Because $\Pr(Q \wedge s) \le \Pr(Q)$ and $\frac{N_{s \wedge q}}{N_s} \le 1$, hence,

$$RF \le \Pr(Q) \log \frac{1}{\Pr(Q)}.$$

Therefore, given an estimate of $\Pr(Q)$, we can find the maximum value of $RF$ from Fig. 2. This gives a range within which a possible threshold can be chosen. We can choose $\delta$ depending on our expectation of the degree of relevance. In general, with a smaller $\delta$, more databases will be considered as relevant. We suggest initially setting $\delta$ as one-third of $RF_{max}$. It can then be adjusted accordingly. In our example, $Q$ is *group = Chinese* and

$\Pr(Q)$ is 40 percent. The corresponding maximum value $RF_{max}$ is about 0.50 from Fig. 2. If we choose $\delta = 0.16$ for both databases, two $RF$ values in the second database are greater than 0.16. According to Definition 1, the second database will be identified as relevant. If we check the values of $RF$ in the first database, none is greater than 0.16. Therefore, it is irrelevant.                □

## 4   ALGORITHM FOR IDENTIFYING RELEVANT DATABASES

With the relevance measure $RF$, we can determine whether a database is relevant to a query predicate before applying data mining algorithms. In this section, we present an algorithm that determines whether a database is relevant to a query predicate $Q$ according to Definition 1.

To compute $RF(s, Q)$, we only need to count three values, $\Pr(Q), \Pr(s)$, and $\Pr(s \wedge Q)$. To determine whether a database is relevant to $Q$, we need to test all selectors, *which can be done by scanning the table once*. Let us assume that there are $m + 1$ attributes, $A_0, \ldots, A_m$. Attribute $A_0$ is referenced by $Q$. For each of the other attributes,
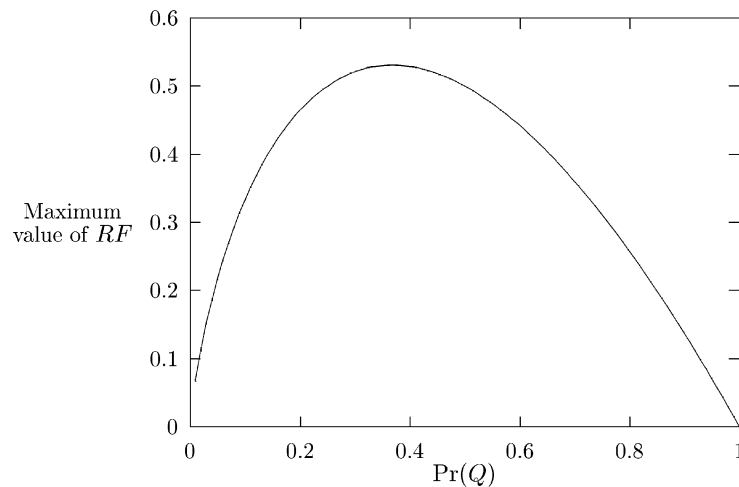


Fig. 2. Maximum value of $RF$ as the function of $\Pr(Q)$.

```
Algorithm RelevantDB


input:       D(A_0, A_1, ..., A_m),          // the database
             Q : A_0 relop Q_value,          // the query predicate
             δ : threshold;                  // the threshold
output:      True or False;                  // TRUE if D is relevant, FALSE otherwise
Q_counter:   number of tuples for which Q holds;


01 begin
02     /* search database to establish counters */
03     foreach record R do begin
04         if R.A_0 = Q_value then
05             increase Q_counter by 1;
06         for attribute A_i, 1 ≤ i ≤ m do begin
07             search record t in Si where t.S_value = R.A_i;
08             if found then
09                 increase t.S_counter by 1;
10             else
11                 insert (R.A_i, 1, 0) into S_i;
12             if R.A_0 = Q_value then
13                 increase t.SQ_counter by 1;
14         end
15     end;
16
17     /* compute RF(s,Q) */
18     foreach table S_i (1 ≤ i ≤ m) do
19         foreach record t do begin
20             RF = ComputeRF(Q_counter, t.S_counter, t.SQ_counter);
21             if RF ≥ δ then
22                 return TRUE;
23         end;
24
25     return FALSE;
26 end;
```

Fig. 3. Algorithm RelevantDB.

$A_i, 1 \le i \le m$, a table $S_i$ is maintained to keep track of selectors and the related counters. An entry of $S_i$ is a triple of $(S\_value, S\_counter, SQ\_counter)$. $S\_value$ is the value of the selector, $S\_counter$ records the number of tuples for which $A_i = S\_value$ is true. $SQ\_counter$ records the number of tuples for which both $Q$ and $A_i = S\_value$ are true. The algorithm is listed in Fig. 3.

Algorithm RelevantDB scans the table and finds all possible selectors and inserts them into the $S$ tables (lines 1-15). Related counters are updated. After the scanning, function $ComputeRF$ is called for each selector to compute its relevance factor, $RF$ (lines 18-20). If $RF$ is greater than the given threshold, the algorithm stops and returns $TRUE$ (lines 21-22). If no selector has $RF$ greater than $δ$, $FALSE$ is returned (line 25), i.e., the database is irrelevant to $Q$. As an example, Fig. 4 lists the contents of the $S$ tables for each attribute after applying algorithm $RelevantDB$ to the first example database in Section 2. From the $S$ tables, relevance factor $RF$ for each selector can be computed.

Algorithm RelevantDB has two parts: 1) Reading each record in the database and 2) searching for the entry of selectors and updating the counters for each attribute in the record. The cost of part 1 is proportional to the number of records in the database. As for part 2, with a proper data structure, for example, using some hashing function for the

selector tables, the search for selector entries can be kept as constant, regardless of the number of selectors of an attribute. Therefore, the run time for the algorithm is $O(NM)$, where $N$ is the number of records in the database and $M$ is the number of attributes.

## 5 EXPERIMENTS

In order to have a better understanding of the issue of relevance and the relevance factor, we conducted two sets of experiments on the relevance measure.

### 5.1 Benchmark Data Sets

There are a large number of benchmark data sets available for evaluating data mining algorithms [22]. However, no similar data sets are available for evaluating measures such as the one proposed in this paper. To obtain multiple possibly relevant data sets, one approach is to vertically partition a data set into a number of data sets, each of which contains a certain number of attributes. It is hoped that some data sets obtained are relevant to a particular query predicate and some are not so that the effectiveness of the proposed measure can be evaluated. In experiments, we use the data sets from the University of California-Irvine repository [2]. The multiple databases are generated from

| alcohol | | |
| --- | --- | --- |
| *S_value* | *S_counter* | *SQ_counter* |
| never | 4 | 2 |
| sometimes | 7 | 3 |
| often | 4 | 1 |

| on_diet | | |
| --- | --- | --- |
| *S_value* | *S_counter* | *SQ_counter* |
| no | 15 | 6 |

| snack-between-meals | | |
| --- | --- | --- |
| *S_value* | *S_counter* | *SQ_counter* |
| seldom | 4 | 2 |
| sometimes | 9 | 3 |
| often | 2 | 1 |

| favorite-non-veg. | | |
| --- | --- | --- |
| *S_value* | *S_counter* | *SQ_counter* |
| pork | 3 | 1 |
| chicken | 3 | 1 |
| beef | 4 | 1 |
| fish | 4 | 2 |
| egg | 1 | 1 |

Fig. 4. Applying RelevantDB to Example DB1.

the Wisconsin breast cancer data and mushroom data. We use this procedure to construct multiple databases:

1. Run a feature ranking algorithm—Relief [14], [5] to rank the attributes of a data set;[1]
2. Divide (vertically partition) the data set into several data sets according to the rankings of the features; and
3. Specify queries according to class labels and see which data sets are identified as relevant.

We should expect that data sets with important features identified by Relief are relevant to some queries regarding the class.

### 5.1.1   Results of the Wisconsin Breast Cancer Data

The data has 10 attributes and the first attribute is ID. Relief ranks importance attributes of the Breast Cancer data in a decreasing order: 7, 3, 4, 8, 6, 2, 5, 9, 10, 1. So, we obtain three data sets: D1 has attributes 7, 3, 4 plus class, D2 has attributes 8, 6, 2 plus class, and D3 has attributes 5, 9, 10, 1 plus class.

The class has two values: "benign," "malignant." In general, people are more concerned about the instances with "malignant" class value. If we specify our query as class = "malignant" and choose $\delta = 0.17$ for three data sets (there are 241 instances which have class label = "malignant," so $\Pr(Q) = 0.34$ and $RF_{max} = 0.51$, thus $\delta = 0.51/3 = 0.17$), we only got one selector: Bare Nuclei (attribute 7) = "10," for which the RF value is 0.277395 from data set 1. Thus, data set D1 was identified as relevant. When we specified the query as class = "benign," three data sets were identified as relevant since selectors with $RF$ larger than 0.17 are obtained from three data sets (Table 6). We ran C4.5rules [27] and confirmed that classification rules for "benign" cases contain attributes 2, 3, 5, 7, and 9 from all three data sets.

### 5.1.2   Results of the Mushroom Data

This data has 22 attributes. Relief ranks importance order of attributes as: 5, 20, 11, 8, 19, 4, 10, 22, 9, 12, 13, 21, 7, 3, 2, 15, 14, 18, 6, 17, 16, 1. We divide the Mushroom data into four data sets: Besides a class attribute, D1 has attributes 5, 20, 11, 8, 19, D2 has attributes 4, 10, 22, 9, 12, 13, D3 has attributes 21, 7, 3, 2, 15, 14, and D4 has attributes 18, 6, 17, 16, 1.

We posed queries based on both class values: "edible," "poisonous" and specified $\delta = 0.16$. The values of $\Pr(Q)$ for "edible," "poisonous" are 0.518 and 0.482 and their $RF_{max}$ values are about 0.48, one third of which is 0.16 as $\delta$. For Query 1: class = "poisonous" and Query 2: class = "edible," we found the selectors shown in Tables 7 and 8, respectively. From the results, it is clear that, for Query 1: class = "poisonous," data sets D1, D2, and D3 are identified as relevant. It should be noted, however, that, for D1 and D2, there are four selectors found for each, but there is only one selector found (population (attribute 21) = "several," RF = 0.19225) in D3. According to the ranking by Relief, attribute 21 is the most important attribute in D3. As for Query 2: class = "edible," data sets D1 and D2 are identified as relevant, D3 and D4 are irrelevant. As we see, the experimental results are consistent with the relevance of the data sets in the two benchmark data sets. We also ran C4.5rules to confirm if any features in D3 and D4 are found in the classification rules. The result is negative: Only features in D1 and D2 are contained in the rules. After the above experiment, we have some understanding of the measure. Now, let us engage another database in further experiment.

### 5.2   A Real-World Data Set

The NSERC (Natural Science and Engineering Research Council of Canada) research grant database is used. It contains eight tables. We briefly explain the largest table, NSERC table. The table contains information on the amount of awards with 14 attributes as follows: *Id, Sysid, Sortname, Dept, Organization-id, Fyr, Compyr, Award, Grant-code, Ctee, Install, Acd3, Discipline-code,* and *Cnt*. To simplify the task of finding meaningful databases, we transform the problem of

---

1. Relief used here is our implementation. It relies on repeated subsampling to weigh each feature. It needs to access data many times and shows good results in removing irrelevant features.

TABLE 6
RF Values for Selectors Found for Query Class = "benign" of the Wisconsin Breast Cancer Data

| D1 | bare nuclei (attribute 7) = 1 | RF=0.307321 |
| | uniformity of cell size (attribute 3) = 1 | RF=0.323374 |
| | uniformity of cell shape (attribute 4) = 1 | RF=0.302165 |
| D2 | single epithelial cell size (attribute 6) = 2 | RF=0.270725 |
| D3 | marginal adhesion (attribute 5) = 1 | RF=0.263845 |
| | normal nucleoli(attribute 9) = 1 | RF=0.270205 |

TABLE 7
RF Values for Selectors Found for Query 1 of the Mushroom Data

| D1 | odor (attribute 5)= "foul" | RF=0.279920 |
| | spore-print-color (attribute 20)= "chocolate" | RF=0.196877 |
| | gill-size (attribute 8)= "narrow" | RF=0.240120 |
| | ring-type(attribute 19)= "large" | RF=0.167952 |
| D2 | bruises (attribute 4)= "no" | RF=0.212519 |
| | gill-color (attribute 9)= "black" | RF=0.223936 |
| | stalk-surface-above-ring (attribute 12)= "silky" | RF=0.263952 |
| | stalk-surface-below-ring (attribute 13)= "silky" | RF=0.255164 |
| D3 | population (attribute 21)= "several" | RF=0.19225 |
| D4 | none | |

TABLE 8
RF Values for Selectors Found for Query 2 of the Mushroom Data

| D1 | odor (attribute 5) = "none" | RF=0.377183 |
| | gill-size (attribute 8)= "broad" | RF=0.208157 |
| | spore-print-color (attribute 20)= "brown" | RF=0.166312 |
| | ring-type (attribute 19)= "pendant" | RF=0.239353 |
| D2 | bruises (attribute 4)= "bruises" | RF=0.221617 |
| | stalk-surface-above-ring(attribute 12)= "smooth" | RF=0.197662 |
| | stalk-surface-below-ring (attribute 13)= "smooth" | RF=0.172113 |
| D3 | none | |
| D4 | none | |

identifying relevant databases from multiple databases to a problem of determining the relevance of a single database to that of different query predicates for a single database.

One of our objectives is to determine whether the table contains interesting information related to the grant amount *Award*. First, we need to determine the threshold. As we discussed in Section 3.3, we estimated the maximum $RF$ as follows: The total number of tuples in the table is 18,510. To make a query predicate have statistical significance, we require that the query cover at least 3 percent of the tuples (about 500). So, we set the lower bound for $\Pr(Q)$ as 0.03. From Fig. 2 we can find that the maximum $RF$ is about 0.15, one third of it is 0.05 as $\delta$.

Since the amount of award is continuous, we construct query predicates in the form of ranges such as $a_1 \leq Award < a_2$. Table 9 lists the testing results of five query predicates. For two queries with the award amount between [0, 20,000) and [20,000, 40,000), we obtain some

TABLE 9
Experiments with NSERC Data

| Query predicate $Q$ | Selector $S$ where $RF(s,Q) > \delta$ | $RF(s,Q)$ |
|---|---|---|
| $0 \leq Award < 20000$ | $Grant\text{-}code =$OGPIN | 0.096328 |
|  | $Sysid =$G | 0.088566 |
|  | $Acd3 =$1201 | 0.060338 |
| $20000 \leq Award < 40000$ | $Grant\text{-}code =$OGPIN | 0.091384 |
|  | $Sysid =$G | 0.091384 |
| $40000 \leq Award < 60000$ | Nil | |
| $60000 \leq Award < 90000$ | Nil | |
| $90000 \leq Award$ | Nil | |

selectors whose $RF$ values are greater than the specified threshold 0.05. For the other three query predicates, no such selector is obtained. In other words, the NSERC table is relevant to the first two queries, but irrelevant to the last three queries. That is, we can find interesting patterns about the award amount in the first two ranges, but not in the last three.

To verify this finding, let us take a closer look at selector $Grant\text{-}code$ = OGPIN (OGPIN refers to an individual research grant type) since the $RF$ value of $Grant\text{-}code$ = OGPIN for both query predicates is greater than the threshold. Fig. 5 shows the distributions of award amount for the top five grant codes: OGPIN, EQPEQ, INFIF, SAPIN, and OGPGP. The area between a curve for a grant code and the x-axis is the number of tuples with that code. As we can see, most tuples with grant code OGPIN have the award amount in the range of (0, 40,000), that is, the distribution of the award amount for OGPIN is special within this range,

therefore, the table is identified as relevant to the queries in these two ranges. Viewing these distributions and their shapes can also determine which code is relevant to the query. However, it would involve human interpretation. It is impractical when the number of codes is large and with many databases to examine. Our measure serves the purpose of human interpretation in an automatic manner.

## 6 RELATED WORK

To our knowledge, little work on mining multiple databases has been reported in the literature. However, research work on the interestingness of discovered knowledge seems quite related to this work, though the objectives and methods are different. As was mentioned earlier, applying data mining techniques to a given database can usually generate some knowledge (or patterns, rules). Typically, the number of such patterns is large and only a small fraction of them may
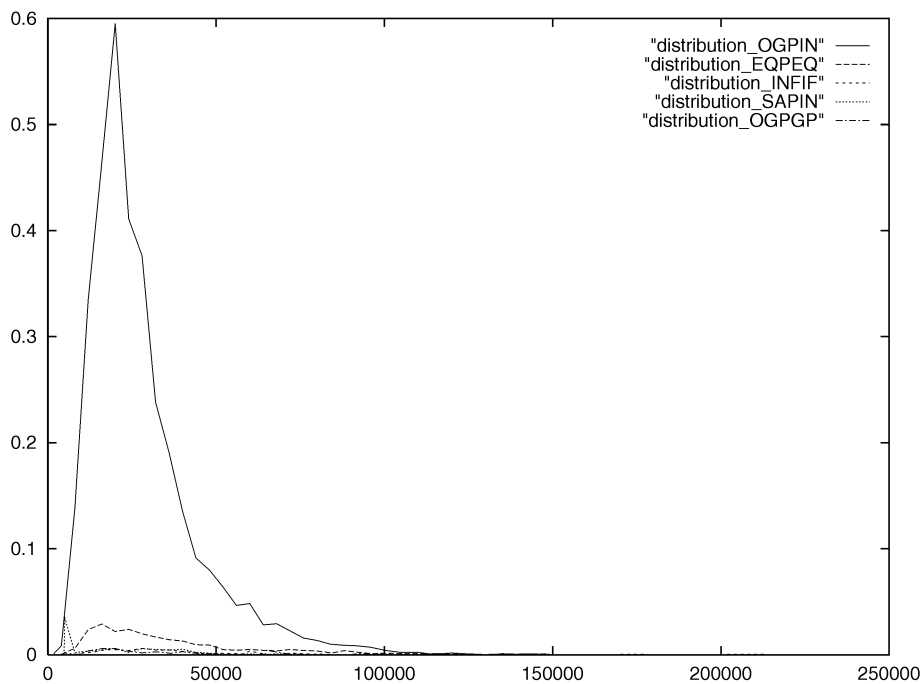


Fig. 5. Distributions of award amount.

TABLE 10
Relevance Factor for Rules in [13]

| Rule | Interesting by $IC^{++}$ | $RF(B_i, A_i)$ |
|---|---|---|
| $A_1 \rightarrow B_1$ | No | 0.00 |
| $A_2 \rightarrow B_2$ | Yes | 0.43 |
| $A_3 \rightarrow B_3$ | Yes | 0.45 |
| $A_4 \rightarrow B_4$ | No | 0.00 |

be of interest to the user. Research work on interestingness of knowledge tries to distinguish the potentially interesting patterns from others. There are a number of approaches studying the interestingness problem. Some researchers proposed approaches that determine interestingness of a discovered rule based on the user's feedback [28], [21], [15], hence, the measurement of interestingness is *subjective*. Some researchers developed various *objective* interestingness measures based on the statistics underlying the discovered patterns [25], [13], [12]. With the quantitative objective interestingness measure, discovered patterns can be ranked and less interesting ones can be filtered out. Kamber and Shinghal [13] gave a good survey on such objective interestingness measures. They also proposed a measure, $IC^{++}$, which is claimed to have the best property regarding evaluating the interestingness of characteristic rules.

While interestingness is an important issue for knowledge discovery, we argue that it is also important to determine whether a database contains potentially interesting information before applying data mining techniques. When we filter out those irrelevant databases, we should not discard the databases that may contain interesting information. To check on this, we apply our relevance measure to the example given in [13]. They presented preliminary results about interestingness of characteristic rules mined from a synthetic database of 4,000 tuples described by 10 attributes. Among four mined rules, $IC^{++}$ indicated that two are interesting and the other two are not interesting. Based on the statistics given, we computed the relevance factors for each rule by taking its lefthand side as the query predicate and the righthand side as the selector. The results are shown in Table 10. It can be seen that our method indeed serves the purpose: If $RF$ is small (0 in the case of Rules 1 and 4), the database is irrelevant and there is no need to mine patterns related to the predicates ($A_1$, $A_4$). Recall that mining rules is more costly than calculating $RF$. If we insist on proceeding, there might be some rules found ($A_1 \rightarrow B_1$ and $A_4 \rightarrow B_4$), but, by applying the interestingness measure $IC^{++}$, they would be considered uninteresting, as shown in the example given by [13]. A major difference between $IC^{++}$ and $RF$ is that the former is applied as a postprocessing filter and the latter as a preprocessing filter.

Another piece of related work can be found in [29] in which the authors developed a cross entropy-based approach to rule induction from databases. It not only learns rules for a given concept (classification), but it simultaneously learns rules relating multiple concepts—known as generalized rule induction. Briefly put, the ITRULE algorithm repeatedly modifies the $K$ rules by calculating their $J$-measure through depth-first search over possible lefthand sides, starting with the first-order conditions. The commonality between our measure and the $J$-measure is that both rely on information computation. Our measure stems directly from the need to cover the five cases in Section 3 in order to select databases instead of rule induction. An important difference is that our measure requires just one pass over a database and no revision is required.

## 7 CONCLUSION

In this paper, we address one issue of data mining from multi-database systems: Identifying those databases relevant to a data mining task from a set of semantically related databases. We argue that effective data mining from multiple databases should involve a preselection phase. In that phase, those databases relevant to the mining task are identified. Data mining techniques are then only applied to those relevant databases. The preselection phase is essential since joining irrelevant databases with relevant ones will enlarge the search space, reduce the statistical significance of knowledge to be discovered, and may lead to too many uninteresting patterns. This work underlines our belief that less can mean more [18].

In order to quantify the relevance of a database, a measurement, $RF$, relevance factor, is proposed. The beauty of $RF$ is that: 1) Its value can indicate the statistical significance of different aspects of an attribute value related to the query predicate and 2) it can be efficiently computed by scanning the database only once. Given also is an algorithm that implements the measurement to check whether a database is relevant to a given query predicate. Experiments were conducted on both benchmark datasets and a real-world database to test the method proposed in the paper and the measure was cross validated by the state-of-the-art interesting measure ($IC^{++}$) in the literature.

Further work will be generalizing the measure to overcome some limitations of the measure. We consider two directions. First, our measure focuses on individual attributes due to the efficiency issue. Therefore, the measure assumes that higher order interdependence among the attributes is somehow reflected in first order attributes. However, further analysis shows the monotonic nature of $RF$ when $\Pr(Q) \leq \Pr(s)$. It may not be so when $\Pr(Q) > \Pr(s)$. Therefore, the above assumption is still required for a general case. Investigation continues. Exhaustive search is infeasible due to its computational complexity. Educated random search may be one way out [24]. Second, moving toward attribute-based identification from selector-based, although we presented a guideline for choosing the threshold value in Section 3, a proper threshold value is also determined by other factors because if query Q covers many selectors of an attribute, $\Pr(Q \wedge s)$ may be much less than $\Pr(Q)$, the upper bound of $RF$ by which we choose the threshold may be relatively too loose. Also, it is possible that one database is identified as relevant just because it contains only one selector with some

high value of $RF$ and another database is filtered out which may contain many selectors with values of $RF$ less than but near the threshold value. In order to prevent this from happening, we plan to introduce a modified measure that sums up $RF$'s whose values are $\epsilon$-within the threshold $\delta$, where $\epsilon$ can be a function of $\delta$, e.g., it's one-third.

One immediate application of this work is for information providers to screen databases they are offered. Imagine that an information dealer for a special trade wants to buy data related to the trade. The dealer pays accordingly for relevant databases and declines irrelevant ones. Our measure $RF$ provides a way for the dealer to screen what to buy or decline.

In addition to the practical value of this work, more opportunities open up for research. One line of research is to exploit irrelevant databases identified by the measure $RF$. We may further probe whether these databases can be merged for reuse. This occurs when we face a particular application with few relevant databases. Obviously, database cleaning and merging requires some methods beyond $RF$. Throughout this paper, we assume that a database for a specific purpose resides in one place. As web technology advances, however, there are cases where one database is distributed in several places. Although it is always possible to consider one database instead of its portions, it is also desirable to deal with each portion. This necessitates an extension to the current form of $RF$ as it is clear that $RF$ cannot directly handle the databases in distributed environments.

## ACKNOWLEDGMENTS

## REFERENCES

[1] R. Agrawal, T. Imielinski, and A. Swami, "Database Mining: A Performance Perspective," *IEEE Trans. Knowledge and Data Eng.,* vol. 5, no. 6, pp. 914-925, Dec. 1993.

[2] C.L. Blake and C.J. Merz, "UCI Repository of Machine Learning Databases," 1998, http://www.ics.uci.edu/~mlearn/MLRepository.html.

[3] L. Breiman, J.H. Friedman, R.A. Olshen, and C.J. Stone, *Classification and Regression.* Wadsworth & Brooks/Cole Advanced & Books Software, 1984.

[4] P. Clark and T. Niblett, "The CN2 Induction Algorithm," *Machine Learning,* vol. 3, pp. 261-283, 1989.

[5] M. Dash and H. Liu, "Feature Selection Methods for Classifications," *Intelligent Data Analysis: An Int'l J.,* vol. 1, no. 3, 1997 (http://www-east.elsevier.com/ida/free.htm).

[6] U.M. Fayyad, G. Piatetsky-Shapiro, and P. Smyth, "From Data Mining to Knowledge Discovery: An Overview," *Advances in Knowledge Discovery and Data Mining,* U.M. Fayyad, G. Piatetsky-Shapiro, P. Smyth, and R. Uthurusamy, eds., pp. 1-34, AAAI Press/The MIT Press, 1996.

[7] U.M. Fayyad, G. Piatetsky-Shapiro, and P. Smyth, "From Data Mining to Knowledge discovery: An Overview," *Advances in Knowledge Discovery and Data Mining,* U.M. Fayyad, G. Piatetsky-Shapiro, P. Smyth, and R. Uthurusamy, eds., pp. 495-515, AAAI Press/The MIT Press, 1996.

[8] J. Han, Y. Cai, and N. Cercone, "Data Driven Discovery of Quantitative Rules in Relational Database," *IEEE Trans. Knowledge and Data Eng.,* vol. 5, no. 1, pp. 29-40, 1993.

[9] J. Han and Y. Fu, "Discovery of Multiple-Level Association Rules from Large Databases," *Proc. 21st VLDB Conf.,* pp. 420-431, 1995.

[10] J. Han and Y. Fu, "Attribute-Oriented Induction in Data Mining," *Advances in Knowledge Discovery and Data Mining,* U.M. Fayyad, G. Piatetsky-Shapiro, P. Smyth, and R. Uthurusamy, eds., pp. 399-421, AAAI Press/The MIT Press, 1996.

[11] J. Hong and C. Mao, "Incremental Discovery of Rules and Structure by Hierarchical and Parellel Clustering," *Knowledge Discovery in Databases,* G. Piatetsky-Shapiro and W.J. Frawley, eds., pp. 177-194, AAAI/The MIT Press, 1991.

[12] F. Hussain, H. Liu, E. Suzuki, and H. Lu, "Exception Rule Mining with a Relative Interestingness Measure," *Proc. Fourth Pacific-Asia Conf. Knowledge Discovery and Data Mining (PAKDD),* 2000.

[13] M. Kamber and R. Shinghal, "Evaluating the Interestingness of Characteristic Rules," *Proc. Second Int'l Conf. Data Mining (KDD-96),* pp. 263-266, 1996.

[14] K. Kira and L.A. Rendell, "The Feature Selection Problem: Traditional Methods and a New Algorithm," *Proc. 10th Nat'l Conf. Artificial Intelligence,* pp. 129-134, 1992.

[15] B. Liu and W. Hsu, "Post-Analysis of Learned Rules," *Proc. 13th Nat'l Conf. Artificial Intelligence (AAAI-96),* pp. 828-834, 1996.

[16] *Feature Extraction, Construction and Selection: A Data Mining Perspective,* H. Liu and H. Motoda, eds. Boston: Kluwer Academic, 1998.

[17] H. Liu and H. Motoda, *Feature Selection for Knowledge Discovery Data Mining.* Boston: Kluwer Academic, 1998.

[18] H. Liu and H. Motoda, "Less Is More," *Feature Extraction, Construction and Selection: A Data Mining Perspective,* H. Liu and H. Motoda, eds. Boston: Kluwer Academic, pp. 3-12, 1998.

[19] H. Liu and R. Setiono, "Feature Selection via Discretization," *IEEE Trans. Knowledge and Data Eng.,* vol. 9, no. 4, pp. 642-645, July/Aug. 1997.

[20] J.A. Major and J. Mangano, "Selecting among Rules Induced from a Hurricane Database," *Proc. AAAI-93 Workshop Knowledge Discovery in Databases,* G. Piatetsky-Shapiro, ed., pp. 28-44, 1993.

[21] C.J. Matheus, G. Piatetsky-Shapiro, and D. McNeill, "Selecting and Reporting What Is Interesting," *Advances in Knowledge Discovery and Data Mining,* U.M. Fayyad, G. Piatetsky-Shapiro, P. Smyth, and R. Uthurusamy, eds., pp. 495-514, AAAI Press/The MIT Press, 1996.

[22] C.J. Merz and P.M. Murphy, "UCI Repository of Machine Learning Databases," http://www.ics.uci.edu/~mlearn/MLRepository.html, Dept. of Information and Computer Science, Univ. of California, Irvine, 1996.

[23] R.S. Michalski, I. Mozetic, J. Hong, and N. Lavrac, "The Multi-Purpose Incremental Learning System aq15 and Its Testing Application to Three Medical Domains," *Proc. Fifth Nat'l Conf. Artificial Intelligence,* pp. 1041-1045, 1986.

[24] R. Motwani and P. Raghavan, "Randomized Algorithms," *The Computer Science and Eng. Handbook,* A.B. Tucker Jr., ed., pp. 141-161, CRC Press and ACM, 1997.

[25] G. Piatetsky-Shapiro, "Discovery, Analysis, and Presentation of Strong Rules," *Knowledge Discovery in Databases,* G. Piatetsky-Shapiro and W.J. Frawley, eds., pp. 229-248, AAAI/The MIT Press, 1991.

[26] G. Piatetsky-Shapiro, C. Matheus, P. Smyth, and R. Uthurusamy, "KDD93: Progress and Challenges," *AI Magazine,* pp. 77-87, Fall 1994.

[27] J.R. Quinlan, *C4.5: Program for Machine Learning.* Morgan Kaufmann, 1993.

[28] A. Silberschatz and A. Tuzhilin, "On Subjective Measures of Interestingness in Knowledge Discovery," *Proc. First Int'l Conf. Knowledge Discovery and Data Mining,* pp. 275-281, 1995.

[29] P. Smyth and R.M. Goodman, "An Information Theorectic Approach to Rule Induction from Databases," *IEEE Trans. Knowledge and Data Eng.,* vol. 4, no. 4, pp. 301-316, Aug. 1992.

[30] R. Srikant and R. Agrawal, "Mining Quantitative Association Rules in Large Relational Tables," *Proc. ACM SIGMOD Conf. Management of Data,* 1996.

[31] R. Uthurusamy, "From Data Mining to Knowledge Discovery: Current Challenges and Future Directions," *Advances in Knowledge Discovery and Data Mining,* U.M. Fayyad, G. Piatetsky-Shapiro, P. Smyth, and R. Uthurusamy, eds., pp. 561-569, AAAI Press/The MIT Press, 1996.

[32] G. Wiederhold, "Foreword: On the Barriers and Future of Knowledge Discovery," *Advances in Knowledge Discovery and Data Mining,* U.M. Fayyad, G. Piatetsky-Shapiro, P. Smyth, and R. Uthurusamy, eds., pp. vii-xi, AAAI Press/The MIT Press, 1996.

[33] R. Zembowicz and J.M. Zytkow, "From Contigency Tables to Various Forms of Knowledge in Databases," *Advances in Knowledge Discovery and Data Mining,* U.M. Fayyad, G. Piatetsky-Shapiro, P. Smyth, and R. Uthurusamy, eds., pp. 329-349, AAAI Press/The MIT Press, 1996.

**Huan Liu** obtained his BEng degree from the Electrical Engineering and Computer Science Department at Shanghai Jiao Tong University in 1983. He earned his PhD degree in computer science in 1989 and master's degree in computer science in 1985 from University of Southern California, Los Angeles. He is currently an associate professor in the Department of Computer Science and Engineering at Arizona State University. Beginning in November 1989, he worked as a researcher/project learder at Telecom Australia Research Laboratories in Melbourne, Australia, for four years. He joined the School of Computing, National University of Singapore, in January 1994, from where he is currently on leave. His major research interests include data reduction and preprocessing, data mining, and real-world applications. He is a senior member of the IEEE and a member of the IEEE Computer Society, ACM, and AAAI. More information can be found at http://www.public.asu.edu/~huanliu.

**Hongjun Lu** received his BSc degree from Tsinghua University, China, and his MSc and PhD degrees from the University of Wisconsin-Madison. He is currently at the Department of Computer Science, Hong Kong University of Science and Technology (on leave from the National University of Singapore). His main research interests are in data/knowledge base management systems with emphasis on query processing and optimization. His recent research work includes data quality, data warehousing, and data mining. He is also interested in the development of Internet-based database applications and electronic business systems. Dr Lu is currently a member of the ACM SIGMOD Advisory Board and ACM SIGKDD International Liaisons. He is the cochair of the steering committee of the Pacific-Asia Conference of Knowledge Discovery and Data Mining and has been serving as a member of the organization and program committees for most important database and KDD conferences, such as ACM SIGMOD, VLDB, ICDE, and ACM SIGKDD.

**Jun Yao** obtained his bachelor's degree in engineering from the Electrical Engineering Department at ZheJiang University in 1987 and his master's degree in computer science from the School of Computingl at National University of Singapore in 1998. He is currently working as a software developer at Mokonet Internet Inc., New York. Before that, he worked as a research fellow/senior research engineer/senior software developer at the National University of Singapore/High Performance Computing Institute, Singapore/Angoss Software Corporation, Canada, respectively.

▷ **For further information on this or any computing topic, please visit our Digital Library at** http://computer.org/publications/dlib.