

Mining Approximative Descriptions of Sets Using Rough Sets

Dan A. Simovici

University of Massachusetts Boston, Dept. of Computer Science,
100 Morrissey Blvd. Boston, Massachusetts, 02125 USA

dsim@cs.umb.edu

Selim Mimaroglu

Bahcesehir University, Dept. of Computer Engineering,
Istanbul, Turkey

selim.mimaroglu@bahcesehir.edu.tr

Abstract

Using concepts from rough set theory we investigate the existence of approximative descriptions of collections of objects that can be extracted from data sets, a problem of interest for biologists that need to find succinct descriptions of families of taxonomic units. Our algorithm is based on an anti-monotonicity of borders of object sets and makes use of an approach that is, in a certain sense, a dual of the Apriori algorithm used in identifying frequent item sets.

1 Introduction

Rough sets are approximative descriptions of sets that can be achieved using equivalences (or partitions). This fertile idea was introduced by the Polish mathematician Z. Pawlak in [7]. Excellent surveys supplemented by large bibliographies are [4] and [1]. The reference [5] contains a rich collection of applications of rough sets.

A set of objects \mathcal{D} is described by a set of attributes K when the set of components of these objects that correspond to K constitute a characteristic of \mathcal{D} . In general, it is difficult to identify sets of objects that are characterizable in this manner and to find the corresponding sets of attributes. If certain errors are tolerable (which is usually the case), then the problem becomes more manageable.

The problem that we address starts with a set of rows of a table (or, in the terminology of rough sets, an information system) and seeks to identify sets of attributes whose values provide approximative descriptions of the set of objects. An area of applications that inspired this work is identifying relevant characters that can describe (even approximatively) a set of biological taxa, where such descriptions serve for specimen identification in the field work. A related problem that involves identifying sets of templates of given length and of minimum support was

discussed in [6]. Similar issues were discussed in the context of formal concept analysis in [3].

The problem and the algorithm are treated using rough set theory, which we introduce in the next sections. After discussing data sets and sets of objects, we introduce our algorithm for mining approximate descriptions. The algorithm is followed by a presentation of preliminary experimental results and we conclude with a discussion of further research topics.

2 Lower and Upper Approximations of Sets and Borders

We begin by introducing some terminology and basic facts from rough set theory. Unless stated otherwise, all sets are finite. If U is a subset of a set S , we denote its complement $S - U$ by U^c .

Let S be a set. An *approximation space* is a pair (S, ρ) , where ρ is an equivalence relation defined on the set S .

DEFINITION 2.1. Let (S, ρ) be an approximation space and let U be a subset of S . The ρ -*lower approximation* of U is the union of all ρ -equivalence classes included in the set U :

$$\text{lap}_\rho(U) = \bigcup \{[x]_\rho \in S/\rho \mid [x]_\rho \subseteq U\}.$$

The ρ -*upper approximation* of U is the union of ρ -equivalence classes that have a nonempty intersection with the set U :

$$\text{uap}_\rho(U) = \bigcup \{[x]_\rho \in S/\rho \mid [x]_\rho \cap U \neq \emptyset\}.$$

□

The following statements hold in an approximation

space (S, ρ) (see, for example [8]):

$$\text{lap}_\rho(\emptyset) = \text{uap}_\rho(\emptyset) = \emptyset, \quad (2.1)$$

$$\text{lap}_\rho(S) = \text{uap}_\rho(S) = S, \quad (2.2)$$

$$\text{lap}_\rho(U \cap V) = \text{lap}_\rho(U) \cap \text{lap}_\rho(V), \quad (2.3)$$

$$\text{uap}_\rho(U \cup V) = \text{uap}_\rho(U) \cup \text{uap}_\rho(V), \quad (2.4)$$

$$\text{lap}_\rho(U \cup V) \supseteq \text{lap}_\rho(U) \cup \text{lap}_\rho(V), \quad (2.5)$$

$$\text{uap}_\rho(U \cap V) \subseteq \text{uap}_\rho(U) \cap \text{uap}_\rho(V), \quad (2.6)$$

$$\text{lap}_\rho(U^c) = (\text{uap}_\rho(U))^c \quad (2.7)$$

$$\text{uap}_\rho(U^c) = (\text{lap}_\rho(U))^c, \quad (2.8)$$

for every $U, V \in \mathcal{P}(S)$.

Note that, in general, $\text{lap}_\rho(U) \subseteq \text{uap}_\rho(U)$ for any set U . Also, we have

$$\text{uap}_\rho(U) = \{t \in S \mid (t, s) \in \rho \text{ for some } s \in U\},$$

$$\text{lap}_\rho(U) = \{t \in U \mid (t, s) \in \rho \text{ implies } s \in U\}.$$

In defining and retrieving approximative descriptions of sets the notion of border plays a fundamental role.

DEFINITION 2.2. The *positive ρ -border* of U is the set

$$\partial_\rho^+(U) = U - \text{lap}_\rho(U),$$

while the *negative ρ -border* of the same set is

$$\partial_\rho^-(U) = \text{uap}_\rho(U) - U,$$

The ρ -border of U is:

$$\partial_\rho(U) = \partial_\rho^+(U) \cup \partial_\rho^-(U) = \text{uap}_\rho(U) - \text{lap}_\rho(U). \quad \square$$

A subset U of S is ρ -rough if $\partial_\rho(U) \neq \emptyset$ and is ρ -crisp otherwise.

Observe that

$$\partial_\rho(U) = \{t \in S \mid (t, s) \in \rho \text{ and } (t, z) \in \rho \text{ for some } s \in U \text{ and } z \notin U\}.$$

For every subset U of an approximation space (S, ρ) we have the equalities:

$$\partial_\rho^+(U^c) = \partial_\rho^-(U), \quad (2.9)$$

$$\partial_\rho^-(U^c) = \partial_\rho^+(U). \quad (2.10)$$

Indeed, we can write

$$\begin{aligned} \partial_\rho^+(U^c) &= U^c - \text{lap}_\rho(U^c) \\ &= U^c - (\text{uap}_\rho(U))^c \\ &\quad (\text{by Equality 2.7}) \\ &= U^c \cap \text{uap}_\rho(U) \\ &= \text{uap}_\rho(U) - U \\ &= \partial_\rho^-(U), \end{aligned}$$

which is Equality 2.9. A similar argument works for Equality 2.10.

Let (S, ρ) be an approximation space and let U and V be two subsets of S . If $U \subseteq V$, then $\text{lap}_\rho(U) \subseteq \text{lap}_\rho(V)$ and $\text{uap}_\rho(U) \subseteq \text{uap}_\rho(V)$.

For our purpose it is especially important to describe the behavior of the lower and upper approximation of a set when the equivalence relations involved are changing.

Let ρ, σ be two equivalence relations on the set S . If $\rho \subseteq \sigma$, then each σ -equivalence class is a ρ -saturated set. Therefore, if $U \subseteq S$ we have

$$\text{lap}_\sigma(U) \subseteq \text{lap}_\rho(U) \subseteq U \subseteq \text{uap}_\rho(U) \subseteq \text{uap}_\sigma(U). \quad (2.11)$$

These inclusions imply

$$\partial_\rho(U) \subseteq \partial_\sigma(U). \quad (2.12)$$

Therefore, we have

$$\partial_{\rho_1 \wedge \rho_2}(U) \subseteq \partial_{\rho_1}(U) \cap \partial_{\rho_2}(U). \quad (2.13)$$

for every subset U of S .

3 Data Sets and Sets of Objects

Let H be a finite set, $H = \{A_1, \dots, A_m\}$. We refer to the elements of H as *attributes* and we assume that for each attribute A_i we have a set that contains at least two elements referred to as the *domain of A_i* and denoted by $\text{Dom}(A_i)$.

A *data set on the set of attributes H* is a function $T : \{1, \dots, n\} \times H \rightarrow \bigcup_{j=1}^m \text{Dom}(A_j)$ such that $T(i, A_j) \in \text{Dom}(A_j)$ for $1 \leq i \leq n$ and $1 \leq j \leq m$.

The sequence $t_k = (T(k, 1), \dots, T(k, m))$ is the k^{th} *object of T* for $1 \leq k \leq n$. The numbers $1, \dots, n$ are the object identifiers (abbreviated as *oids*). The set of objects of T is the set $\mathcal{O}_T = \{t_1, \dots, t_n\}$.

If $L = \{A_{i_1}, \dots, A_{i_p}\}$ be a subset of H . The projection of the object $t_k = (T(k, 1), \dots, T(k, m))$ on L is the p -tuple $(T(k, i_1), \dots, T(k, i_p))$.

DEFINITION 3.1. Let T be a data set and let L be a set of attributes of T . The equivalence ρ_L on \mathcal{O}_T defined by

$$\rho_L = \{(t, t') \in \mathcal{O}_T^2 \mid t[L] = t'[L]\}. \quad \square$$

It is easy to see that if L and K are two attribute sets then $\rho_{KL} = \rho_K \cap \rho_L$, where KL is an alternative notation for the union $K \cup L$ of the two attribute sets. Therefore, if $L \subseteq K$, then $\rho_K \subseteq \rho_L$, so by the inclusion (2.12) we have:

$$\partial_{\rho_K}(U) \subseteq \partial_{\rho_L}(U). \quad (3.14)$$

In other words, the border of a set of objects relative to an attribute set is anti-monotonic with respect to the attribute

set. To simplify notations we denote the set $\partial_{\rho_K}(U)$ by $\partial_K(U)$.

DEFINITION 3.2. A set of objects \mathcal{D} is *described by a set of attributes* K if $\partial_K(\mathcal{D}) = \emptyset$ and we refer to K as an *exact description* of \mathcal{D} . \square

The notion of *template* is defined as a formula $\tau = (A_{i_1} = a_{i_1}) \wedge (A_{i_2} = a_{i_2}) \wedge \dots \wedge (A_{i_p} = a_{i_p})$, where $h \neq k$ implies $A_{i_h} \neq A_{i_k}$ for $1 \leq h, k \leq p$ (see, for example [6]). The *length* of the template is p and the *support* of the template τ in the data set T is

$$\text{supp}_T(\tau) = |\{t \in \mathcal{O}_T \mid t[A_{i_1} \dots A_{i_p}] = (a_{i_1}, \dots, a_{i_p})\}|.$$

Note that each template $\tau = (A_{i_1} = a_{i_1}) \wedge (A_{i_2} = a_{i_2}) \wedge \dots \wedge (A_{i_p} = a_{i_p})$ that has a non-zero support corresponds to an equivalence class of ρ_K , where $K = \{A_{i_1} \dots A_{i_p}\}$.

It is shown in [6], starting from the problem of the balanced complete bipartite subgraph [2] that, for a bipartite graph $G = (V_1 \cup V_2, E)$ and two positive integers $k_1 \leq |V_1|$ and $k_2 \leq |V_2|$, it is an NP-complete problem to determine the existence of two subsets U_1 and U_2 of V_1 and V_2 , respectively such that $|U_1| = k_1$, $|U_2| \geq k_2$ and $\{u_1, u_2\} \in E$ for any $u_1 \in U_1$ and $u_2 \in U_2$. This variant of the problem is known as the Complete Bipartite Subgraph (CBS) problem. Then, given a table and two positive integers k and ℓ , the existence of a template τ containing ℓ conjunctions and having support at least k is polynomially equivalent to the CBS problem and therefore, is NP-complete.

We focus in this paper on approximate descriptions of sets of objects which can be stated as follows: given a data set T , a set of objects $\mathcal{D} \subseteq \mathcal{O}_T$, we seek to determine whether there exists an attribute set K containing no more than k attributes such that the size of border $\partial_K(\mathcal{D})$ is less than p . Thus, a search for a template of minimum support is replaced by the search for a set of attributes of limited size that describes the set \mathcal{D} with a certain precision.

Example. Let T be a data set having the set of attributes $H = ABCD$. Our search space is the set of subsets of H . Suppose that we seek to identify a description of the set $\mathcal{D} = \{t_5, t_6, t_7, t_8, t_9\}$, where T is shown in Table 1. Observe that the equivalence classes of the equivalence ρ_{BC} are $\{t_1, t_2\}$, $\{t_3, t_4\}$, $\{t_5, t_6\}$, $\{t_7, t_8\}$, $\{t_9, t_{10}\}$, $\{t_{11}\}$, and $\{t_{12}\}$. Two of these classes, $\{t_5, t_6\}$, $\{t_7, t_8\}$ are included in \mathcal{D} and therefore, they constitute the lower approximation of \mathcal{D} . The class $\{t_9, t_{10}\}$ intersects both \mathcal{D} and its complement and therefore, $\partial_{BC}(\mathcal{D}) = \{t_9, t_{10}\}$. Also, it is clear that $\partial_{BC}^+(\mathcal{D}) = \{t_9\}$ and $\partial_{BC}^-(\mathcal{D}) = \{t_{10}\}$.

Table 1. Data Set T and the set of objects

$$\mathcal{D} = \{t_5, t_6, t_7, t_8, t_9\}$$

	T			
t_1	a_1	b_2	c_1	d_1
t_2	a_2	b_2	c_1	d_2
t_3	a_3	b_1	c_2	d_1
t_4	a_4	b_1	c_2	d_3
t_5	a_1	b_1	c_1	d_2
t_6	a_3	b_1	c_1	d_2
t_7	a_5	b_3	c_3	d_4
t_8	a_1	b_3	c_3	d_2
t_9	a_2	b_3	c_2	d_3
t_{10}	a_3	b_3	c_2	d_3
t_{11}	a_4	b_2	c_2	d_1
t_{12}	a_1	b_3	c_4	d_4

4 An Algorithm for Mining for Approximate Descriptions

In practice approximative descriptions are sufficient and we give an algorithm that identifies such descriptions. This algorithm is, in a certain sense, a dual of the well-known Apriori-algorithm for finding frequent item sets.

DEFINITION 4.1. Let ϵ be a number such that $0 \leq \epsilon \leq 1$. A set of objects \mathcal{D} is ϵ -described by a set of attributes K if

$$\frac{|\partial_K(\mathcal{D})|}{|\mathcal{D}|} \leq \epsilon.$$

\square

The “dual” of the Rymon tree used in the study of frequent item sets is introduced next.

DEFINITION 4.2. Let $H = \{A_1, \dots, A_n\}$ be a set of attributes. The *dual Rymon tree* of H is a rooted tree having $\mathcal{P}(H)$ as its set of nodes, H as its root, and whose set of edges consist of pairs of the form $(U, V) \in \mathcal{P}(H)$ such that V is obtained by dropping from U an attribute that follows the attributes of $H - U$ in the list (A_1, \dots, A_n) . \square

Example. The dual Rymon tree of the set of attributes $H = \{A, B, C, D\}$ is shown in Figure 1.

Our algorithm makes use of the dual Rymon tree of the sets of attributes of the data set and is based on the anti-monotonicity of the size of the border of a set given by the inclusion (3.14). The goal of the algorithm is to compute the smallest sets of attributes that yield a description of a set of objects whose error is below a prescribed threshold.

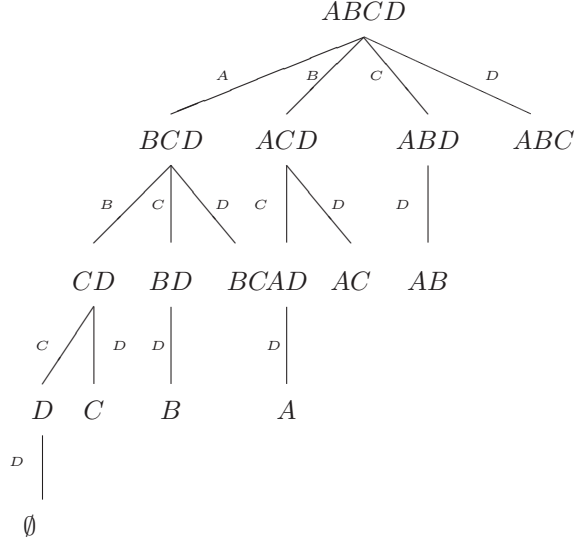


Figure 1. Rymon tree of the set of attributes
 $H = \{A, B, C, D\}$

For identifying approximative descriptions of \mathcal{D} we search the Rymon tree of the set of attributes in a breath-first manner. During the search we apply a pruning technique that reduces the size of the search space substantially.

Starting at the root node of the tree, computation of negative and positive borders take place in breadth first search fashion. In a database having no duplicates the error of the root node is zero. The error of the children of the root are computed next. If the error of a node K is greater than the error threshold there is no need for computing the negative and positive borders for its descendants (which are subsets of the set of the node) because of the anti-monotonicity property of the size of the border contained by Inequality 3.14. Thus, we can prune all descendants of K .

First, a technique for computing the border set of a set of objects $\mathcal{D} = \{t_{i_1}, \dots, t_{i_p}\}$ of a data set T determined by a set of attributes K is introduced in Figure 2. Let $\mathcal{D}^c = \mathcal{O}_T - \mathcal{D} = \{t_{j_1}, \dots, t_{j_q}\}$ be the set of objects of T that do not belong to \mathcal{D} . The projection of each object t_{i_k} in \mathcal{D} on the set of attributes K is compared to the same projection of each object t_{j_h} in \mathcal{D}^c as shown in the Figure 2.

Pruning is implemented as shown in Figure 3. Failed set of attributes are stored according to their cardinality. Before adding any set of attributes into the working queue for computing the negative and positive borders, we check the failed set of attributes. A set of attributes L is not added to the queue if it has a superset that failed.

The main algorithm is shown in Figure 4

Input: T : data set, \mathcal{D} : set of objects, K : set of attributes

Output: Positive (Pos) and Negative (Neg) borders of \mathcal{D}

```

1  $Pos := \{\}$ ;
2  $Neg := \{\}$ ;
3  $\mathcal{D}^c := \mathcal{O}_T - \mathcal{D}$ ;
4 foreach  $t \in \mathcal{D}$  do
5   foreach  $t' \in \mathcal{D}^c$  do
6     // project on  $K$ 
7     if  $t[K] == t'[K]$  then
8       add  $t$  to  $Pos$ ;
9       add  $t'$  to  $Neg$ ;
9 output  $Pos \cup Neg$ ;

```

Figure 2. Computation of borders,
 $FindBorder(T, \mathcal{D}, K)$

Input: L : list of failed descriptors, R : set of attributes

Output: all the qualified $|R| - 1$ size children of R

```

1 enumerate all unique  $|R| - 1$  size children of  $R$  into  $P$ ;
2 foreach  $p \in P$  do
3   if  $L$  contains a superset of  $p$  then
4     remove  $p$  from  $P$ ;
5 output  $P$ ;

```

Figure 3. Pruning of attribute sets,
 $Pruning(L, R)$

5 Experimental Results

Our algorithm is applicable to data sets that have attributes with finite domains of arbitrary cardinality. However, we focus on data sets with binary domains in order to take advantage of the efficient use of memory that bit vectors allow and of the speed of bit operations. If an attribute A has a non-binary domain, $\text{Dom}(A) = \{a_1, \dots, a_k\}$, we replace A by k attributes A^1, \dots, A^k and we replace the A -component of an object t by k components $t[A^1], \dots, t[A^k]$ defined by

$$t[A^j] = \begin{cases} 1 & \text{if } t[A] = a_j, \\ 0 & \text{otherwise} \end{cases}$$

for $1 \leq j \leq k$. Thus, an object of a binary data set is a sequence $t \in \{0, 1\}^n$.

If we define the characteristic n -tuple $r_K =$

Input: T : data set, H : set of attributes, \mathcal{D} : set of objects, err : error threshold

Output: all the descriptors of \mathcal{D}

```

1 initialize a queue  $Q$ ;
2 initialize a list  $L$ ;
3 add  $H$  to  $Q$ ;
4 while  $Q$  is not empty do
5    $R :=$  remove first element from  $Q$ ;
6   if  $FindBorder(T, \mathcal{D}, R) \leq err$  then
7     output  $R$ ;
8      $children := Pruning(L, R)$ ;
9     add  $children$  to  $Q$ ;
10  else
11    add  $R$  to  $L$ ;
```

Figure 4. Find descriptors of \mathcal{D} , $FindAll(T, H, \mathcal{D}, err)$

$(r_1, \dots, r_n) \in \{0, 1\}^n$ of a set of attributes K as

$$r_i = \begin{cases} 1 & \text{if } A_i \in K, \\ 0 & \text{otherwise.} \end{cases}$$

Thus, the condition $t[K] == t'[K]$ can be expressed as $(t \wedge r_K) \oplus (t' \wedge r_K) = \mathbf{0}$, where $\mathbf{0} = (0, \dots, 0)$. If $t[K]$ equals $t'[K]$, then t is added to the positive border $\partial_K^+(D)$ and t' is added to the negative border $\partial_K^-(D)$. After computing the positive and the negative borders, we obtain the border $\partial_K(D) = \partial_K^+(D) \cup \partial_K^-(D)$.

To evaluate the effectiveness of our algorithm we conducted a set of preliminary experiments on synthetically generated binary databases using a Pentium 3.0GHz computer having 4GB of main memory running on Linux. We implemented our algorithm in Java which offers support for bit vectors and bit vector operations and we use the advantage of the speed of bit operations.

Figures 9-12 show the running time results on various size databases for a defining set that contains 100 objects. Experiments show that for even big size database run times for reasonable error rates are practical. This is due to our pruning technique and use of bit vectors. It is interesting to note (see Figure 5– 8) that keeping constant the size of the collection of objects that is to be characterized ($|\mathcal{D}| = 100$), for the same error level, the number of descriptors is decreasing when the size of the data set is increasing because of the dramatic increase in the size of the negative border.

In the same time, the larger negative border increases the size of the collection of “failed” descriptors and leads to the apparently paradoxical behavior points to a decrease in computation time for larger data sets (see Figures 9– 12).

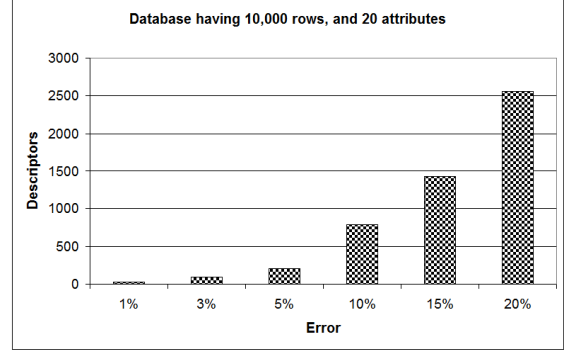


Figure 5. Unique descriptors for a defining set of size 100

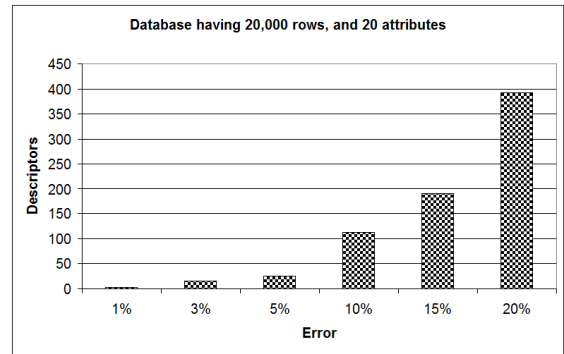


Figure 6. Unique descriptors for a defining set of size 100

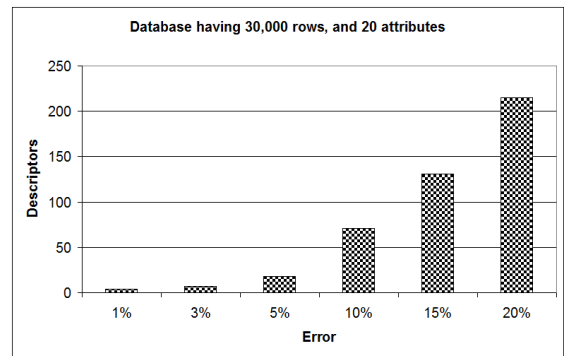


Figure 7. Unique descriptors for a defining set of size 100

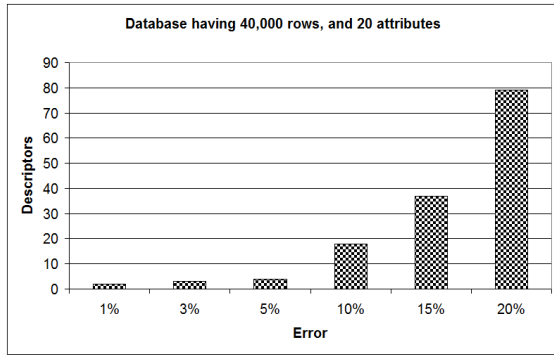


Figure 8. Unique descriptors for a defining set of size 100

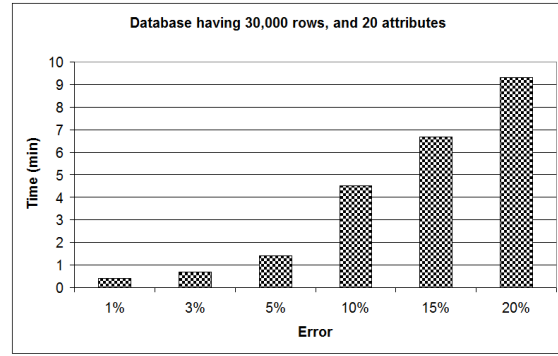


Figure 11. Run-time results for a defining set of size 100

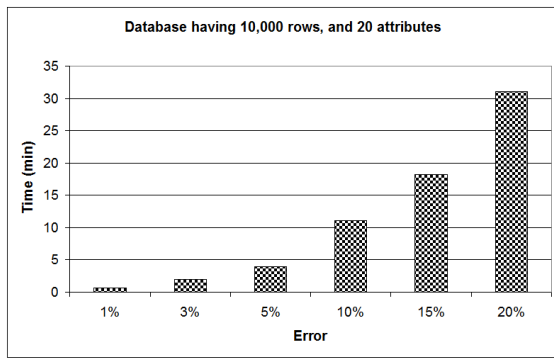


Figure 9. Run-time results for a defining set of size 100

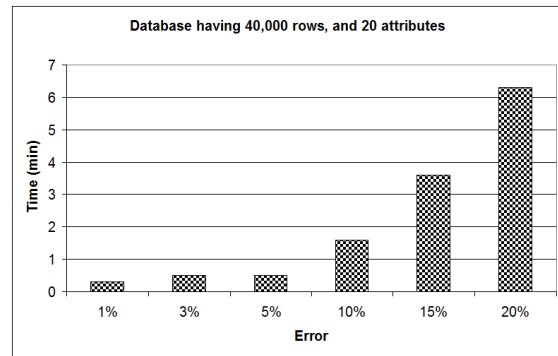


Figure 12. Run-time results for a defining set of size 100

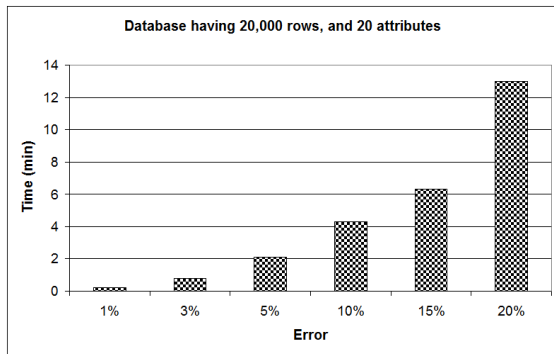


Figure 10. Run-time results for a defining set of size 100

6 Conclusions and Future Work

This paper introduces an efficient algorithm that determines sets of attributes whose values offer approximative descriptions of given sets of objects within a certain margin of

error. It is an open problem to prove that finding exact descriptions of such sets of attributes is NP-complete.

We intend to run further experimental work involving binarization of large taxonomic database and extracting approximative descriptions of large sets of taxa. It should be interesting to contrast approximative descriptions obtained by imposing separate limitations on the positive border and the negative border because such limitations may be important for applications where the importance of false positives or false negatives may be vastly different.

Introducing costs for using specific attributes is also a topic worth investigating since these costs may very considerably and may be, in limited cases, prohibitive.

References

- [1] I. Düntsch and G. Gediga. *Rough Sets Analysis - A Road to Non-invasive Knowledge Discovery*. Methodos, Bangor, UK, 2000.
- [2] M. R. Garey and D. S. Johnson. *Computers and Intractabil-*

- ity – *A Guide to the Theory of NP-Completeness*. W.H. Freeman, New York, 1979.
- [3] B. Ganter and S. O. Kuznetsov. Scale coarsening as feature selection. In *Formal Concept Analysis*, volume LNAI 4933, pages 217–228. Springer-Verlag, Berlin, 2008.
 - [4] J. Komorowski, Z. Pawlak, L. Polkowski, and A. Skowron. Rough sets: A tutorial. In S. K. Pal and A. Skowron, editors, *Rough Sets Hybridization: A New Trend in Decision Making*, pages 3–98. Springer-Verlag Telos, 1999.
 - [5] T. Y. Lin and N. Cercone, editors. *Rough Sets and Data Mining*. Kluwer Academic Publishers, Boston, 1997.
 - [6] S. H. Nguyen, A. Skowron, and P. Synak. Discovery of data patterns with applications to decomposition and classification problems. In *Rough Sets in Knowledge Discovery*, volume 2, pages 55–97. Physica-Verlag, Heidelberg, 1998.
 - [7] Z. Pawlak. *Rough Sets: Theoretical Aspects of Reasoning About Data*. Kluwer Academic Publishing, Dordrecht, 1991.
 - [8] D. Simovici and C. Djeraba. *Mathematical Tools for Data Mining*. Springer-Verlag, London, UK, 2008.