

RSES and RSESLib - A Collection of Tools for Rough Set Computations

Jan G. Bazan¹ and Marcin Szczuka²

¹ Institute of Mathematics, Pedagogical University of Rzeszów
Rejtana 16A, 35-310 Rzeszów, Poland
e-mail: bazan@univ.rzeszow.pl

² Institute of Mathematics, Warsaw University
Banacha Str. 2, 02-097, Warsaw, Poland
e-mail: szczuka@mimuw.edu.pl

Abstract. Rough Set Exploration System - a set of software tools featuring a library of methods and a graphical user interface is presented. Methods, features and abilities of the implemented software are discussed and illustrated with a case study in data analysis.

1 Introduction

Research in decision support systems, classification algorithms in particular those concerned with application of rough sets requires experimental verification. At certain point it is no longer possible to perform every single experiment using software designed for a single purpose. To be able to make thorough, multi-directional practical investigations one have to possess an inventory of software tools that automatise basic operations, so it is possible to focus on the most essential matters. That was the idea behind creation of Rough Set Exploration System, further referred as RSES for short.

First version of RSES and the library RSESLib was released several years ago. The RSESLib is also used in computational kernel of ROSETTA - an advanced system for data analysis (see [16]) constructed at NTNU (Norway) which contributed a lot to RSES development and gained wide recognition. Comparison with other classification systems (see [11]) proves its value.

The RSES software and its computational kernel - the new RSESLib 2.0 library maintains all advantages of previous version. The algorithms from the first incarnation of library are now re-mastered to provide better flexibility, extended functionality and ability to process massive data sets. New algorithms added to the library reflect the current state of our research in classification methods originating in rough sets theory.

The library of functions is not sufficient as an answer to experimenters' demand for helpful tool. Therefore the RSES user interface was constructed. This interface allows to use RSESLib interactively.

2 Basic notions

In order to provide clear description further in the paper we bring here some essential definitions

Information system ([12]) is a pair of the form $\mathbf{A} = (U, A)$ where U is a *universe of objects* and $A = (a_1, \dots, a_m)$ is a set of *attributes* i.e. mappings of the form $a_i : U \rightarrow V_a$, where V_a is called *value set* of the attribute a_i . The decision table is also a pair of the form $\mathbf{A} = (U, A \cup \{d\})$ where the major feature that is different from the information system is the distinguished attribute d . We will further assume that the set of decision values is finite. The i -th *decision class* is a set of objects $C_i = \{o \in U : d(o) = d_i\}$, where d_i is the i -th decision value taken from decision value set $V_d = \{d_1, \dots, d_{rank(d)}\}$

For any subset of attributes $B \subset A$ *indiscernibility relation* $IND(B)$ is defined as follows:

$$xIND(B)y \Leftrightarrow \forall_{a \in B} a(x) = a(y) \quad (1)$$

where $x, y \in U$.

Having indiscernibility relation we may define the notion of reduct. $B \subset A$ is a *reduct* of information system if $IND(B) = IND(A)$ and no proper subset of B has this property. In case of decision tables the *decision reduct* is a set $B \subset A$ of attributes such that it cannot be further reduced and $IND(B) \subset IND(d)$.

Decision rule is a formula of the form

$$(a_{i_1} = v_1) \wedge \dots \wedge (a_{i_k} = v_k) \Rightarrow d = v_d \quad (2)$$

where $1 \leq i_1 < \dots < i_k \leq m$, $v_i \in V_{a_i}$. Atomic subformulae $(a_{i_1} = v_1)$ are called *conditions*. We say that rule r is *applicable* to object, or alternatively, the object *matches* rule, if its attribute values satisfy the premise of the rule. *Support* denoted as $Supp_{\mathbf{A}}(r)$ is equal to the number of objects from \mathbf{A} for which rule r applies correctly $Match_{\mathbf{A}}(r)$ is the number of objects in \mathbf{A} for which rule r applies in general. Analogously the notion of matching set for a rule or collection of rules may be introduced (see [2], [4]).

By *cut* for an attribute $a_i \in A$, such that V_{a_i} is an ordered set we will denote a value $c \in V_{a_i}$. With the use of cut we may replace original attribute a_i with new, binary attribute which tells as whether actual attribute value for an object is greater or lower than c (more in [8]).

Template of \mathbf{A} is a propositional formula $\bigwedge (a_i = v_i)$ where $a_i \in A$ and $v_i \in V_{a_i}$. A generalised template is the formula of the form $\bigwedge (a_i \in T_i)$ where $T_i \subset V_{a_i}$. An object *satisfies* (matches) a template if for every attribute a_i occurring in the template the value of this attribute on considered object is equal to v_i (belongs to T_i in case of generalised template). The template induces in natural way the split of original information system into two distinct subtables containing objects that do or do not satisfy the template, respectively. Decomposition tree is a binary tree, whose every internal node is labeled by some template and external node (leaf) is associated with a set of objects matching all templates in a path from the root to a given leaf (see [7]).

3 The RSES library v. 2.0

The RSES library (RSESLib) is constructed according to the principles of object oriented programming. The programming language used in implementation is Microsoft Visual C++ compliant with ANSI/ISO C++ (ISO/IEC 14882 standard).

The algorithms that have been implemented in the RSES library fall into two general categories.

First category gathers the algorithms aimed at management and edition of data structures that are present in the library.

The algorithms for performing Rough Set theory based operations on data constitute the second, most essential kind of tools implemented inside RSES library. To give the idea what apparatus is given to the user we describe shortly the most important algorithms.

Reduction algorithms i.e algorithms allowing calculation of the collections of reducts for a given information system (decision table). The exhaustive algorithm for calculation of all reducts is present, however such operation may be time-consuming due to computational complexity of such task (see [13]). Therefore approximate and heuristic solutions such as genetic or Jonhson algorithms were implemented (see [14], [6] for details). The library methods for calculation of reducts allow setting initial conditions for number of reducts to be calculated, required accuracy, coverage and so on. Basing on calculated reduct it is possible to calculate decision rules (see [4]). Procedures for rule calculation allow user to determine some crucial constrains for the set of decision rules. Rules received are accompanied with several coefficients that are further used while the rules are being applied to the set of objects (see [3], [2]). In connection with algorithms for reduct/rule calculation appear the subclass of algorithms allowing shortening of rules and reducts with respect to different requirements (see [3]).

Discretisation algorithms allow to find cuts for attributes. In this way initial decision table is converted to one described with less complex, binary attribute without lose of information about discernibility of objects (see [10], [8], [2]).

Template generation algorithms provide means for calculation of templates and generalised templates. Placed side by side with template generation are the procedures for inducing table decomposition trees (see [7] and [9]).

Classification algorithms used for establishing decision value with use of decision rules and/or templates. Operations for voting among rules with use of different schemes fall into this category (see [3], [9], [4], [2]).

During operation certain functions belonging to RSESLib may read and write information to/from files. Most of the files that can be read or written are regular ASCII text files. They particular sub-types can be distinguished by reviewing the contents or identifying file extensions.

4 The RSES GUI.

To simplify the use of RSES library and make it more intuitive a graphical user interface was constructed. This interface allows interaction with library methods

in two modes. First is directed towards ease of use and visual representation of workflow. The other is intended to provide tools for construction of scripts that use RSES functionality

4.1 The project interface.

Project interface window (see Figure 1) consists of two parts. Upper part is the project workspace where icons representing objects occurring during our computation are presented. Lower part is dedicated to messages, status reports, errors and warnings produced during operations. It was designers intention to

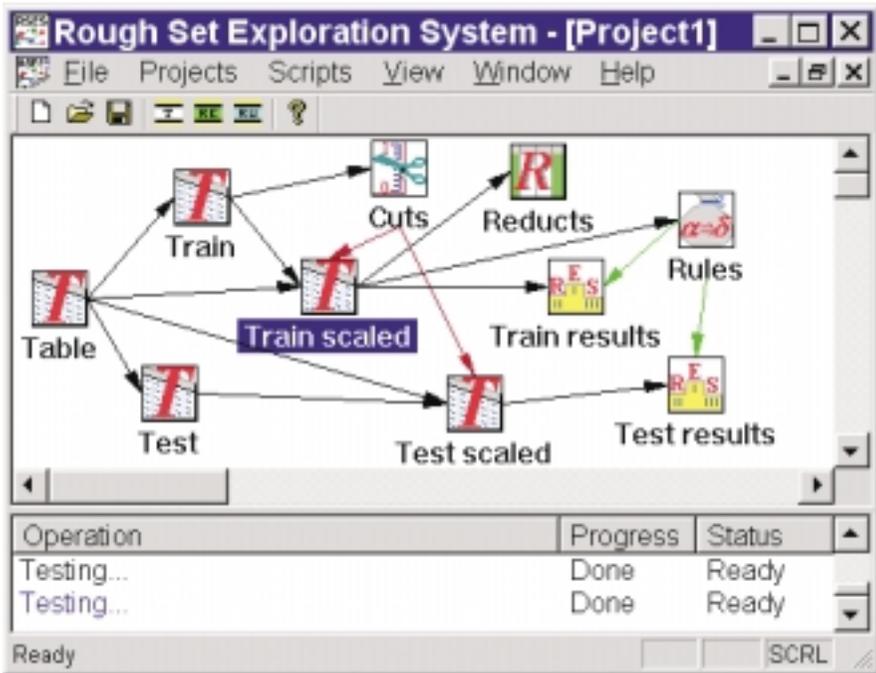


Fig. 1. The project interface window.

simplify the operations on data within project. Therefore, the entities appearing in the process of rough set based computation are represented in the form of icons placed in the upper part of workplace. Such an icon is created every time the data (table, reducts, rules,...) is loaded from the file. User can also place an empty object in the workplace and further fill it with results of operation performed on other objects. The objects that may exist in the workplace are: decision table, collection of reducts, set of rules, decomposition tree, set of cuts and collection of results. Every object (icon) appearing in the project have a set of actions connected with it. By right-clicking on the object the user invokes

a context menu for that object. It is also possible to call the necessary action from general pull-down program menu in the main window. Menu choices allow to view and edit objects as well as make them input to some computation. In many cases choice of some option from context menu will cause a new dialog box to open. In this dialog box user can set values of coefficients used in desired calculation, in particular, designate the variable which will store the results of invoked operation. If the operation performed on the object leads to creation of new object or modification of existing one then such a new object is connected with edge originating in object (or objects) which contributed to its current state. Setting of arrows connecting icons in the workspace changes dynamically as new operations are being performed.

The entire project can be saved to file on disk to preserve results and information about current setting of coefficients. That also allows to re-create the entire work on other computer or with other data.

4.2 Scripting interface.

In case we have to perform many experiments with different parameter settings and changing data it is more convenient to plan such a set of operations in advance and then let computer calculate. The idea of simplifying the preparation and execution of compound experiments drove the creation of RSES scripting mechanism.

The mechanism for performing script-based operations with use of RSES library components consists of three major parts. The scripting interface is the part visible to user during script preparation. Other two are behind the scenes and perform simple syntax checking and script execution. We will not describe checking and executing in greater detail.

The user interface for writing RSES based scripts is quite simple. Main window is split into two parts of which upper contains workplace where scripts are being edited and lower contains messages generated during script execution (Figure 2).

The RSES scripting language constructs available to user are:

- Variables. Any variable is inserted to script with predefined type. The type may be either standard (e. g. integer, real) or RSES-specific.
- Functions. User can use all the functions from RSES armory as well as standard arithmetic operations. Function in script always returns a value.
- Procedures. The procedures from RSES library, unlike functions, may not return a value. The procedures correspond to major operations such as: reduct calculation, rule shortening, loading and saving objects.
- Conditional expressions. The expressions of the form **If ... then ... else** may be used. The user defines condition with use of standard operations.
- Loops. Simple loop can be used within RSES script. The user is required to designate loop control variable and set looping parameters.

While preparing a script the user may not freely edit it. He/she can only insert or remove one of the constructs mentioned above using context menu which

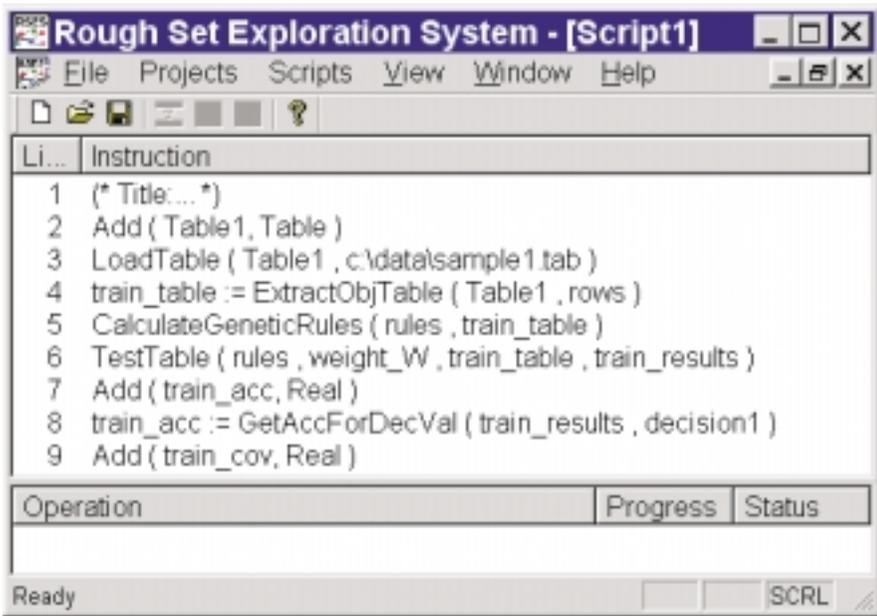


Fig. 2. The script interface window.

appears after right-clicking on the line in the script where insertion/removal is due to occur. The new construct may also be inserted with use of toolbar. The operation of inserting the new construct involves choosing the required operation name and setting all required values. All this is done by easy point-and-click operations within appropriate dialog box supported by pull-down lists of available names. The edition operations are monitored by checker to avoid obvious errors.

Once the script edition is finished it can be executed with menu command **Run**. Before execution the syntax is checked once again. The behaviour of currently running script may be seen in the lower part of interface window.

5 Case study - decomposition by template tree

As already mentioned, the ability of dealing with large data sets is one of key new features of RSESLib 2.0. To deal with such a massive data we use decomposition based on template-induced trees.

Decomposition is a problem of partitioning a large data table into smaller ones. One can use templates extracted from data to partition data table into blocks of objects with common features. We consider here decomposition schemes based on a template tree (see [7]). The main goal of this method is to construct a decomposition tree. Let \mathbf{A} be a decision table. The algorithm for decomposition tree construction can be presented as follows:

Algorithm 1 *Decomposition by template tree* (see [7])

Step 1 Find the best template T in \mathbf{A} .

Step 2 Divide \mathbf{A} onto two subtables: \mathbf{A}_1 containing all objects satisfying T and $\mathbf{A}_2 = \mathbf{A} - \mathbf{A}_1$.

Step 3 If obtained subtables are of acceptable size (in the sense of rough set methods) then stop
else repeat 1-3 for all "too large" subtables.

This algorithm produces a binary tree of subtables with corresponding sets of decision rules for subtables in the leaves of the tree.

The decision tree produced by algorithm presented below can be used to classify a new case to proper decision class. Suppose we have a binary decomposition tree. Let u by a new object and $\mathbf{A}(T)$ be a subtable containing all objects matching template T . We classify object u starting from the root of the tree as follows:

Algorithm 2 *Classification by template tree* (see [7])

Step 1 If u matches template T found for \mathbf{A}
then: go to subtree related to $\mathbf{A}(T)$
else: go to subtree related to $\mathbf{A}(\neg T)$.

Step 2 If u is at the leaf of the tree then go to 3
else: repeat 1-2 substituting $\mathbf{A}(T)$ (or $\mathbf{A}(\neg T)$) for \mathbf{A} .

Step 3 Classify u using decision rules for subtable attached to the leaf

This algorithm uses a binary decision tree, however it should not be mistaken for C4.5, ID3. As we told before, in our experiments (see Section 5.1), a rough set methods have been used for classifying algorithm construction in leaves of the decomposition tree (see [2] and [3] for more details).

5.1 Experiments with Forest CoverType data

The Forest CoverType data used in our experiments were obtained from US Forest Service (USFS) Region 2 Resource Information System (RIS) data (see [5] [17]). Data is in rectangle form with 581012 rows and 56 columns (55 + decision). There are 7 decision values. This data has been studied before using Neural Networks and Discriminant Analysis Methods (see [5]).

The original Forest CoverType data was divided into a training set (11340 objects), a validation set (3780 objects) and a testing set (565892 objects) (see [5]). In our experiments we used algorithms presented in Section 5 that are implemented in RSEslib. The decomposition tree for data have been created only by reference to the training set. The validation set was used for adaptation of classifying algorithms which obtained in the leaves of the decomposition tree (see [2] and [3] for more). The classification algorithm was applied to new cases from the test set. As a measure of classification success we use *accuracy* (see e.g. [7], [5]). The accuracy we define as the ratio of the number of properly classified new cases to the total number of new cases. Three different classification systems applied to the Forest CoverType data have given accuracy of 0.70 (Neural Network), 0.58 (Discriminant Analysis) and 0.73 (RSES).

Acknowledgement: In the first place the special tribute should be paid to Professor Andrzej Skowron who, for all these years, is the *spiritus movens* of RSES evolution. We want to stress our gratitude to colleagues who greatly contributed providing their expertise: Adam Cykier, Nguyen Sinh Hoa, Nguyen Hung Son, Jakub Wróblewski, Piotr Synak, Aleksander Ørn.

Development of our software was significantly supported by ESPRIT project 20288, KBN grant 8T11C02511 and Wallenberg Foundation - WITAS project.

References

1. Skowron A., Polkowski L.(ed.), Rough Sets in Knowledge Discovery 1 & 2, Physica Verlag, Heidelberg, 1998
2. Bazan J., A Comparison of Dynamic and non-Dynamic Rough Set Methods for Extracting Laws from Decision Tables, In [1], Vol. 1, pp. 321-365
3. Bazan J., Approximate reasoning methods for synthesis of decision algorithms (in Polish), Ph. D. Thesis, Department of Math., Comp. Sci. and Mechanics, Warsaw University, Warsaw, 1998
4. Bazan J., Son H. Nguyen, Trung T. Nguyen, Skowron A. and J. Stepaniuk, Decision rules synthesis for object classification. In: E. Orłowska (ed.), Incomplete Information: Rough Set Analysis, Physica - Verlag, Heidelberg, 1998, pp. 23-57.
5. Blackard, J.,A., Comparison of Neural Networks and Discriminant Analysis in Predicting Forest Cover Types, Ph. D. Thesis, Department of Forest Sciences, Colorado State University. Ford Collins, Colorado, 1998.
6. Garey M., Johnson D., Computers and Intractability: A Guide to the Theory of NP-completeness, W.H. Freeman&Co., San Francisco, 1998, (twentieth print)
7. Nguyen Sinh Hoa, Data regularity analysis and applications in data mining. Ph. D. Thesis, Department of Math., Comp. Sci. and Mechanics, Warsaw University, Warsaw, 1999
8. Nguyen Sinh Hoa, Nguyen Hung Son, Discretization Methods in Data Mining, In [1], Vol. 1, pp. 451-482
9. Hoa S. Nguyen, A. Skowron and P. Synak, Discovery of data patterns with applications to decomposition and classification problems. In [1], Vol. 2, pp. 55-97.
10. Nguyen Hung Son, Discretization of real value attributes. Boolean reasoning approach. Ph. D. Thesis, Department of Math., Comp. Sci. and Mechanics, Warsaw University, Warsaw, 1997
11. Michie D., Spiegelhalter D. J., Taylor C. C., Machine Learning, Neural and Statistical Classification, Ellis Horwood, London, 1994
12. Pawlak Z., Rough Sets: Theoretical Aspects of Reasoning about Data, Kluwer, Dordrecht, 1991
13. Rauszer C., Skowron A., The Discernibility Matrices and Functions in Information Systems, In: Słowiński R. (ed.), Intelligent Decision Support, Kluwer, Dordrecht 1992.
14. Wróblewski J., Covering with Reducts - A Fast Algorithm for Rule Generation, Proceeding of RSCTC'98, LNAI 1424, Springer Verlag, Berlin, 1998, pp. 402-407
15. Bazan J., Szczuka M., The RSES Homepage, <http://alfa.mimuw.edu.pl/~rses>
16. Ørn A., The ROSETTA Homepage, <http://www.idi.ntnu.no/~aleks/rosetta>
17. Bay, S. D. , The UCI KDD Archive, <http://kdd.ics.uci.edu>