



UNIVERSITI PUTRA MALAYSIA

**NEW LEARNING MODELS FOR GENERATING CLASSIFICATION
RULES BASED ON ROUGH SET APPROACH**

LUAI ABDEL LATEEF AL SHALABI

FSKTM 2000 2

**NEW LEARNING MODELS FOR GENERATING CLASSIFICATION
RULES BASED ON ROUGH SET APPROACH**

By

LUAI ABDEL LATEEF AL SHALABI

**Thesis Submitted in Fulfilment of the Requirement for the
Degree of Doctor of Philosophy in the
Faculty of Computer Science and Information Technology
Universiti Putra Malaysia**

September 2000



**Dedicated to my father; Abdel Lateef,
my mother; Faidah,
my wife; Samah and the family**

Abstract of thesis presented to the Senate of Universiti Putra Malaysia in fulfilment of the requirements of the degree of Doctor of Philosophy

**NEW LEARNING MODELS FOR GENERATING CLASSIFICATION
RULES BASED ON ROUGH SET APPROACH**

By

LUAI ABDEL LATEEF AL SHALABI

September 2000

Chairman: Ramlan Mahmud, Ph.D.

Faculty: Computer Science and Information Technology

Data sets, static or dynamic, are very important and useful for presenting real life features in different aspects of industry, medicine, economy, and others. Recently, different models were used to generate knowledge from vague and uncertain data sets such as induction decision tree, neural network, fuzzy logic, genetic algorithm, rough set theory, and others. All of these models take long time to learn for a huge and dynamic data set. Thus, the challenge is how to develop an efficient model that can decrease the learning time without affecting the quality of the generated classification rules. Huge information systems or data sets usually have some missing values due to unavailable data that affect the quality of the generated classification rules. Missing values lead to the difficulty of extracting useful information from that data set. Another challenge is how to solve the problem of missing data.

Rough set theory is a new mathematical tool to deal with vagueness and uncertainty. It is a useful approach for uncovering classificatory knowledge and building a classification rules. So, the application of the theory as part of the learning models was proposed in this thesis.

Two different models for learning in data sets were proposed based on two different reduction algorithms. The split-condition-merge-reduct algorithm (SCMR) was performed on three different modules: partitioning the data set vertically into subsets, applying rough set concepts of reduction to each subset, and merging the reducts of all subsets to form the best reduct. The enhanced-split-condition-merge-reduct algorithm (ESCMR) was performed on the above three modules followed by another module that applies the rough set reduction concept again to the reduct generated by SCMR in order to generate the best reduct, which plays the same role as if all attributes in this subset existed. Classification rules were generated based on the best reduct.

For the problem of missing data, a new approach was proposed based on data partitioning and function mode. In this new approach, the data set was partitioned horizontally into different subsets. All objects in each subset of data were described by only one classification value. The mode function was applied to each subset of data that has missing values in order to find the most frequently occurring value in each attribute. Missing values in that attribute were replaced by the mode value.

The proposed approach for missing values produced better results compared to other approaches. Also, the proposed models for learning in data sets generated the

classification rules faster than other methods. The accuracy of the classification rules by the proposed models was high compared to other models.

Abstrak tesis yang dikemukakan kepada Senat Universiti Putra Malaysia sebagai memenuhi keperluan untuk ijazah Doktor Falsafah

**MODEL PEMBELAJARAN BARU UNTUK MENJANA
PETUAPENGGKELASAN BERDASARKAN PENDEKATAN SET KASAR**

Daripada

LUAI ABDEL LATEEF AL AHALABI

September 2000

Pengerusi: Ramlan Mahmud, Ph.D

Fakulti: Sains Komputer dan Teknologi Maklumat

Set-set data, statik atau dinamik, adalah sangat penting dan berguna untuk memaparkan fitur realiti hidup dalam berlainan aspek di bidang industri, perubatan, ekonomi dan lain-lain. Kebelakangan ini, pelbagai model telah digunakan untuk menjana pengetahuan daripada set-set data yang kabur dan tidak pasti seperti pokok keputusan induksi, rangkaian neural, logik kabur, algoritma genetik, teori set kasar dan sebagainya. Kesemua model tersebut mengambil masa yang panjang untuk belajar bagi set-set data yang besar dan dinamik. Maka, cabarannya ialah bagaimana untuk membangunkan satu model yang berkesan yang boleh mengurangkan masa pembelajarannya tetapi tidak menjejaskan kualiti petua klasifikasi yang terjana. Maklumat yang besar atau set-set data yang besar selalu mempunyai nilai yang tercicir akibat daripada ketidaksempurnaan data yang menyebabkan kerumitan untuk mengekstrak maklumat yang berfaedah daripada set-set data itu.

Teori set kasar adalah alat matematik baru untuk mengatasi masalah kekaburan dan ketidakpastian. Ia adalah pendekatan yang berguna untuk mendapatkan pengetahuan pengkelasan dan membina satu petua pengkelasan. Malah, aplikasi teori set kasar sebagai sebahagian daripada model pembelajaran telah diperkenalkan dalam tesis ini.

Dua model pembelajaran berlainan dalam set-set data telah dicadangkan berdasarkan kepada dua algoritma reduksi berlainan. Algoritma SCMR telah dilakukan ke atas tiga modul berlainan: pembahagian set-set data kepada subset, aplikasikan konsep reduksi set kasar kepada setiap subset, dan gabungan hasil reduksi semua subset untuk membentuk reduksi terbaik. Model ESCMR telah digunakan ke atas tiga modul di atas dan diikuti dengan satu modul lain yang mengaplikasikan sekali lagi konsep reduksi set kasar ke atas reduksi yang terhasil daripada SCMR, untuk menjana reduksi terbaik. Reduksi ini memainkan peranan yang sama seperti semua atribut di dalam subset. Petua pengkelasan dijana berdasarkan reduksi terbaik.

Pendekatan baru telah dicadangkan untuk mengatasi masalah ini. Dalam pendekatan baru ini, set-set data telah dikategorikan secara mendatar kepada subset yang berlainan. Semua objek dalam setiap subset data ditakrifkan oleh hanya satu nilai klasifikasi. Fungsi mod telah diaplikasikan ke atas setiap subset data yang mengandungi nilai tercicir untuk mencari nilai yang mempunyai frekuensi tertinggi dalam setiap atribut. Nilai tercicir dalam atribut itu digantikan dengan nilai mod.

Pendekatan yang dicadangkan untuk nilai tercacir ini menghasilkan keputusan yang lebih baik berbanding dengan pendekatan lain. Model yang dicadangkan dalam pembelajaran set-set data, telah dapat menjana peraturan klasifikasi lebih cepat daripada kaedah lain. Ketepatan petua pengkelasan oleh model yang dicadangkan adalah lebih tinggi berbanding dengan model-model yang lain.

ACKNOWLEDGEMENTS

Praise to Allah for giving me the strength and patience to complete this moderate work.

Firstly, I would like to thank my supervisor Dr. Ramlan Mahmod, deputy dean, Faculty of Computer Science and Information Technology, for his supervision and invaluable discussion. My thanks also to my co-supervisors, Dr. Abdul Azim Abd ghani, dean, Faculty of Computer Science and Information Technology and Dr. Md. Yazid Mohd. Saman for their invaluable discussion and help.

I am certainly grateful to Prof. Ohmann, Theoretical Surgery Unit, Department of General and Trauma Surgery, Heinrich-Heine-University, Dusseldorf, Germany, for providing the original Acute Abdominal Pain data set.

I also wish to express my thanks to the Faculty of Computer Science and Information Technology, Post Graduate Office, and Library, University Putra Malaysia, for providing assistance and good environment while I was pursuing my research work.

Special thanks are due to my friends and colleagues who have made some assistance in the development of this research. From list I would like to mention my friends Ayman Oklat and Ziad Abu Kaddorah.

Finally, I would like to thank my parents, wife, sisters, and brother for their patience, understanding and moral support.

Luai Abdel Lateef Al Shalabi

September, 2000

TABLE OF CONTENTS

		Page
DEDICATION.....		ii
ABSTRACT.....		iii
ABSTRAK.....		vi
ACKNOWLEDGEMENTS.....		ix
APPROVAL SHEETS.....		xi
DECLARATION.....		xiii
LIST OF TABLES.....		xviii
LIST OF FIGURES.....		xix
LIST OF ABBREVIATIONS.....		xx
CHAPTER		
1	INTRODUCTION.....	1.1
	1.1 Background.....	1.1
	1.2 Objectives of the Study.....	1.2
	1.3 Significance of the Study.....	1.3
	1.4 Contribution of the Study.....	1.5
	1.5 Applications.....	1.6
	1.6 Organisation of the Thesis.....	1.6
2	LITERATURE REVIEW.....	2.1
	2.1 Introduction.....	2.1
	2.2 Data Mining.....	2.3
	2.3 Aspects of Data Mining.....	2.3
	2.3.1 Uncertainty Handling.....	2.4
	2.3.2 Efficiency of Algorithms.....	2.4
	2.3.3 Constraining Knowledge Discovered.....	2.4
	2.3.4 Incorporating Domain Knowledge.....	2.5
	2.3.5 Size and Complexity of Data.....	2.5
	2.3.6 Data Selection.....	2.5
	2.3.7 Understandability of Discovered Knowledge.....	2.6
	2.3.8 Consistency Between Data and Discovered Knowledge.....	2.6
	2.4 Training set and testing set.....	2.6
	2.5 Learning the Knowledge.....	2.7
	2.5.1 Deductive and Inductive Learning Methods.....	2.7
	2.5.2 Supervised and Unsupervised Learning.....	2.8
	2.6 Representation of Knowledge.....	2.9
	2.7 Rule Generation.....	2.12
	2.7.1 Semantic of Definite Rules.....	2.13
	2.7.2 Non-Monotonic and Default Reasoning.....	2.13
	2.8 Classification of Data Mining Problems.....	2.15
	2.8.1 Classification.....	2.15

2.8.2	Association.....	2.17
2.8.3	Sequences.....	2.17
2.9	Data Mining Model	2.18
2.9.1	Data Pre-Processing.....	2.18
2.9.2	Data Mining Tools.....	2.18
2.9.3	User Bias.....	2.19
2.10	Problems and Challenges in Data Mining.....	2.19
2.10.1	Noisy Data.....	2.19
2.10.1.1	Corrupted Values.....	2.19
2.10.1.2	Missing Attributes Values.....	2.20
2.10.2	Difficult Training Set.....	2.20
2.10.2.1	Non-Representative Data.....	2.21
2.10.2.2	Absence of Boundary Cases.....	2.21
2.10.2.3	Limited Information.....	2.21
2.10.3	Dynamic Data Sets.....	2.22
2.10.4	Huge Data Sets.....	2.22
2.11	Data Mining Methods.....	2.23
2.11.1	Neural Networks Approach.....	2.23
2.11.2	Bayesian Approach.....	2.24
2.11.3	Decision Tree Approach.....	2.25
2.11.4	Rough Set Approach.....	2.27
2.11.4.1	ROUGHIDAS.....	2.27
2.11.4.2	ROUGHCLASS.....	2.28
2.11.4.3	ROUGHIAN.....	2.29
2.11.5	Other Methods.....	2.30
3	ROUGH SETS.....	3.1
3.1	Introduction.....	3.1
3.2	Rough Set Theory.....	3.3
3.3	Problems Addressed by Rough Set Theory.....	3.5
3.4	Concepts of Rough Set Theory.....	3.6
3.4.1	Information System.....	3.6
3.4.2	Indiscernibility Relation.....	3.7
3.4.3	Approximation of Sets.....	3.9
3.4.4	Rough Classification.....	3.11
3.4.5	Reduction and Dependency of Attributes.....	3.12
3.4.6	Decision Tables and Decision Rules.....	3.14
4	DESIGN OF LEARNING SYSTEM USING ROUGH SETS.....	4.1
4.1	Introduction.....	4.1
4.2	Problems of Missing Data and the New Approach.....	4.2
4.3	Testing of the New Approach.....	4.5
4.4	Preparing the Data Set for the Proposed Models	4.9
4.4.1	Acquisition of Data.....	4.9
4.4.2	Formalization of Decision Tables.....	4.9
4.4.3	Reduction of Attributes.....	4.10
4.5	Models for Learning Classification Rules in Data Sets.....	4.14
4.6	Strategy of the Proposed Models.....	4.19
4.7	Application to Hepatitis Data	4.20
4.8	Evaluation of Rules and Techniques	4.22

4.9	Advantages of the Proposed Models	4.24
5	EXPERIMENTS AND RESULTS.....	5.1
5.1	Introduction.....	5.1
5.2	Data for Application.....	5.2
5.3	Organisation of Data.....	5.3
5.4	System Learning.....	5.5
5.5	Experiment 1	5.7
5.5.1	Training The Three Models Using the Original Information System.....	5.8
5.5.1.1	Model Based on Rough Sets.....	5.8
5.5.1.2	Model Based on SCMR.....	5.9
5.5.1.3	Model Based on ESCMR.....	5.10
5.5.2	Observations.....	5.10
5.5.3	Training The Three Models After the Deletion of Age Attribute.....	5.11
5.5.3.1	Model Based on Rough Sets.....	5.11
5.5.3.2	Model Based on SCMR.....	5.12
5.5.3.3	Model Based on ESCMR.....	5.13
5.5.4	Observations.....	5.13
5.5.5	Training The Three Models After the Deletion of Spongiosis Attribute.....	5.14
5.5.5.1	Model Based on Rough Sets.....	5.14
5.5.5.2	Model Based on SCMR.....	5.15
5.5.5.3	Model Based on ESCMR.....	5.15
5.5.6	Observations.....	5.16
5.6	Final Remarks	5.17
5.7	Experiment 2.....	5.28
5.7.1	Training The Three Models Using the Original Information System.....	5.28
5.7.1.1	Model Based on Rough Sets.....	5.28
5.7.1.2	Model Based on SCMR.....	5.29
5.7.1.3	Model Based on ESCMR.....	5.29
5.7.2	Observations.....	5.30
5.7.3	Training The Three Models After the Deletion of Age Attribute.....	5.31
5.7.3.1	Model Based on Rough Sets.....	5.31
5.7.3.2	Model Based on SCMR.....	5.31
5.7.3.3	Model Based on ESCMR.....	5.32
5.7.4	Observations.....	5.32
5.7.5	Training The Three Models After the Deletion of Spongiosis Attribute.....	5.33
5.7.5.1	Model Based on Rough Sets.....	5.34
5.7.5.2	Model Based on SCMR.....	5.34
5.7.5.3	Model Based on ESCMR.....	5.35
5.7.6	Observations.....	5.35
5.7.7	Final Remarks.....	5.36
5.8	Experiment 3	5.46
5.8.1	Data Collection.....	5.47
5.8.2	Observations.....	5.47

5.8.3	Final Remarks	5.51
6	CONCLUSION AND FUTURE WORK.....	6.1
6.1	Conclusion.....	6.1
6.2	Capabilities of the Models.....	6.4
6.3	Future Work.....	6.5
	BIBLIOGRAPHY.....	R.1
	VITA.....	V.1

LIST OF TABLES

Table	Page
3.1	An information system of credit card applications..... 3.7
3.2	Equivalence classes generated from IND(C)..... 3.8
3.3	The lower approximation of sets Y1 and Y2 of accepted and rejected applications..... 3.11
3.4	Reduced decision table (R1) for credit card applications..... 3.16
3.5	Reduced decision table (R2) for credit card applications..... 3.16
4.1	Evaluation of the proposed approach using small data set..... 4.6
4.2	Evaluation of the proposed approach using hepatitis data set..... 4.7
4.3	Division of medical test data of hepatitis..... 4.21
4.4	A brief comparison of 5 different models (only for Class=Live)..... 4.22
5.1	Division of medical test data of dermatology..... 5.4
5.2	Division of medical test data of acute abdominal pain..... 5.5
5.3	Results of experiment 1 using all the attributes..... 5.11
5.4	Results of experiment 1 after the deletion of the Age attribute..... 5.14
5.5	Results of experiment 1 after the deletion of the Spongiosis attribute..... 5.16
5.6	Time spent to construct the best reduct after the deletion of the Age attribute..... 5.18
5.7	Time spent to construct the best reduct after the deletion of the Spongiosis attribute..... 5.18
5.8	Number of classification rules discovered using the models..... 5.19
5.9	Accuracy of the models..... 2.20
5.10	Results of experiment 2 using all the attributes..... 5.30
5.11	Results of experiment 2 after the deletion of the Age attribute..... 5.33
5.12	Results of experiment 2 after the deletion of the Spongiosis attribute..... 5.36
5.13	Time spent to construct the best reduct after the deletion of the Age attribute..... 5.37
5.14	Time spent to construct the best reduct after the deletion of the Spongiosis attribute..... 5.37
5.15	Number of classification rules discovered using the models..... 5.37
5.16	Overall accuracy of the learning systems..... 5.49
5.17	Disease specific accuracy of the learning systems..... 5.50
5.18	Distribution of the diagnostic categories for the training set..... 5.53

LIST OF FIGURES

Figure	Page
3.1	The approximation space for accepted credit card applications..... 3.11
4.1	New approach for missing values..... 4.4
4.2	Accuracy comparison for the small data set..... 4.6
4.3	Accuracy comparison for the huge data set..... 4.8
4.4	A model based on SCMR 4.16
4.5	A model based on enhanced SCMR..... 4.17
5.1	Construction time (using all attributes in the data set)..... 5.21
5.2	Construction time (using all attributes except Age)..... 5.21
5.3	Construction time (using all attributes except Spongiosis)..... 5.22
5.4	Number of attributes (using all attributes in the data set)..... 5.22
5.5	Number of attributes (using all attributes except Age)..... 5.23
5.6	Number of attributes (using all attributes except Spongiosis)..... 5.23
5.7	Number of rules (using all attributes in the data set)..... 5.24
5.8	Number of rules (using all attributes except Age)..... 5.24
5.9	Number of rules (using all attributes except Spongiosis)..... 5.25
5.10	Time spent after updating the data set by removing the Age attribute..... 5.26
5.11	Time spent after updating the data set by removing the Spongiosis attribute..... 5.27
5.12	Construction time (using all attributes in the data set)..... 5.39
5.13	Construction time (using all attributes except Age)..... 5.39
5.14	Construction time (using all attributes except Spongiosis)..... 5.40
5.15	Number of attributes (using all attributes in the data set)..... 5.40
5.16	Number of attributes (using all attributes except Age)..... 5.41
5.17	Number of attributes (using all attributes except Spongiosis)..... 5.41
5.18	Number of rules (using all attributes in the data set)..... 5.42
5.19	Number of rules (using all attributes except Age)..... 5.42
5.20	Number of rules (using all attributes except Spongiosis)..... 5.43
5.21	Time spent after updating the data set by removing the Age attribute..... 5.44
5.22	Time spent after updating the data set by removing the Spongiosis attribute..... 5.45

LIST OF ABBREVIATIONS

AI	Artificial Intelligence
BR	Best Reduct
ESCMR	Enhanced Split Condition Merge Reducts
FOL	First Order Logic
ID3	Induction of Decision Trees
IS	Information System
ITRule	Information Theoretic Rule Induction
PRISM	Programming In Statistical Modeling
RoughDAS	Rough sets based Data Analysis System
RoughIAN	Rough Information Analysis
RSDA	Rough Set Data Analysis
SCMR	Split Condition Merge Reduct
TDIDT	Top-Down Induction of Decision Tree

CHAPTER 1

INTRODUCTION

1.1 Background

Data mining or knowledge discovery from databases is a form of machine discovery where the discovered knowledge is represented in a high level language. It is capable of discovering domain knowledge from given examples. However, the theory of knowledge discovery is still under development and few existing methods are practical. The type of rule or pattern that exists in data depends on the domain. Discovery systems have been applied to real databases in medicine, astronomy, the stock market, and many others.

Data mining has come to refer to the process of analysing data and generating new knowledge that is previously hidden and unseen. The overall goal is to create a simplified model of the domain under study. Various techniques for data mining have been employed, mostly from the area of inductive learning, with different forms of knowledge representation such as weights in artificial neural network or nodes in a decision tree.

A well-known and widely employed inductive learning algorithm is ID3 (Quinlan, 1979; 1986). ID3 and subsequent versions (Quinlan, 1993) based their decisions on statistical measures of the information entropy and knowledge is represented as a decision tree which may easily be converted to a set of rules. An alternative method is based on rough sets which is based on the theory of sets and

topology. Rough set theory was introduced by Pawlak (Pawlak, 1982) and since then a number of applications have been reported in diverse fields such as medical diagnosis, conflict analysis, and process control (Slowinski, 1992). Rough set can be used for the purpose of generating *If...Then* rules and/or as a technique for eliminating redundant information prior to the use of, say, artificial neural networks.

Rough set theory, introduced by Pawlak in 1982 (Pawlak, 1991; 1995), is a new mathematical tool dealing with vagueness and uncertainty. It has proved its soundness and usefulness in many real life applications. Rough set theory offers effective methods that are applicable in many branches of AI. The idea of rough set consists on approximation of a set by a pair of sets called lower and upper approximations of the set. The definition of the approximation follows from an indiscernibility relation between elements of the sets, called objects. Objects are described by attributes of a qualitative nature.

Rough set can be a useful tool for pre-processing data for generating classification rules. Applying its concepts to an application at hand reduces the number of attributes and the complexity of the classification rules (Al-Shalabi et al., 1999b).

1.2 Objectives of the Study

The main goal of this dissertation is to propose new models based on time and cost for learning classification rules from data sets. Solving the problem of repeating a learning process to the whole original data set if it changes is a special case of the proposed models. This learning process takes some time to generate new

classification rules. As we know most data sets are frequently changed. This kind of data sets is called “dynamic” data sets. Time consuming is the important disadvantages of the existing systems that are used to generate classification rules especially if the data set is huge. In this study, medical diagnosis has been applied.

The objectives of the study can be derived from the main goal, and they include:

1. To produce a new approach for pre-processing input with missing data to the proposed models.
2. To produce new algorithms for reducing a number of attributes of a data set based on the rough set theory and the partitioning of the data set.
3. To produce new models for the discovery of accurate classification rules from data sets based on the new algorithms.

1.3 Significance of the Study

Data sets, static or dynamic, are very important and useful for presenting real life features in different aspects of industry, medicine, economy, and others. They act as store of information that can answer most questions. In order to gain the most knowledge from these data sets, a special technique to sieve and clean these data is used and then correct knowledge is generated. Building an expert system is one of these techniques. It is considered one of the most important and early used systems

that can generate knowledge from databases. Recently, many other techniques have been used to generate knowledge from databases or data sets such as ID3, neural network, fuzzy logic, genetic algorithm, rough set theory, and others. We may say that all these methods have the capability of generating knowledge requested from databases or data sets. But these methods seem to have learning time and learning complexity problems especially if the data set is dynamic. So, this raises a question of how to develop an efficient tool that can help the task and decrease the learning time without affecting the quality of the rules generated.

This study proposes new different models to approach this problem. These have the idea of partitioning the data set vertically into two or more subsets. Then a rough set theory is applied to each subset in the new space in order to find the best reduct for each of them. The best reduct contains less number of attributes, strong discernibly attributes, low cost, and some other features. The combination of all reducts is done in the higher level and as a result the knowledge presented by classification rules are discovered.

The concept of rough sets has been proposed as a new mathematical tool to deal with uncertain and imprecise data, and it seems to be of significant importance to AI and cognitive sciences (Slowinski, 1992). Using this tool to approach the problem of data reduction and data dependency has emerged as a powerful technique in the application of expert systems, decision support systems, machine learning, and pattern recognition.