



UNIVERSITI PUTRA MALAYSIA

**ATTRIBUTE SET WEIGHTING AND DECOMPOSITION
APPROACHES FOR REDUCT COMPUTATION**

QASEM AHMAD AL-RADAIDEH.

FSKTM 2005 7

**ATTRIBUTE SET WEIGHTING AND DECOMPOSITION APPROACHES
FOR REDUCT COMPUTATION**

By

QASEM AHMAD AL-RADAIDEH

**Thesis Submitted to the School of Graduate Studies, Universiti Putra Malaysia,
in Fulfilment of the Requirements for the Degree of Doctor of Philosophy**

July 2005



DEDICATION

*To the memory of my Father,
To my great Mother,
To my Wife and Daughters,
To my Brother and Sisters.*

Qasem



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

**ATTRIBUTE SET WEIGHTING AND DECOMPOSITION APPROACHES
FOR REDUCT COMPUTATION**

By

QASEM AHMAD AL-RADAIDEH

July 2005

Chairman : Associate Professor Md. Nasir Sulaiman, PhD

Faculty : Computer Science and Information Technology

This research is mainly in the Rough Set theory based knowledge reduction for data classification within the data mining framework. To facilitate the Rough Set based classification, two main knowledge reduction models are proposed. The first model is an approximate approach for object reducts computation used particularly for the data classification purposes. This approach emphasizes on assigning weights for each attribute in the attributes set. The weights give indication for the importance of an attribute to be considered in the reduct. This proposed approach is named Object Reduct by Attribute Weighting (ORAW). A variation of this approach is proposed to compute full reduct and named Full Reduct by Attribute Weighting (FRAW).

The second proposed approach deals with large datasets particularly with large number of attributes. This approach utilizes the principle of incremental attribute set decomposition to generate an approximate reduct to represent the entire dataset. This proposed approach is termed for Reduct by Attribute Set Decomposition (RASD).



The proposed reduct computation approaches are extensively experimented and evaluated. The evaluation is mainly in two folds: first is to evaluate the proposed approaches as Rough Set based methods where the classification accuracy is used as an evaluation measure. The well known *10-fold* cross validation method is used to estimate the classification accuracy. The second fold is to evaluate the approaches as knowledge reduction methods where the size of the reduct is used as a reduction measure.

The approaches are compared to other reduct computation methods and to other none Rough Set based classification methods. The proposed approaches are applied to various standard domains datasets from the UCI repository. The results of the experiments showed a very good performance for the proposed approaches as classification methods and as knowledge reduction methods. The accuracy of the ORAW approach outperformed the Johnson approach over all the datasets. It also produces better accuracy over the Exhaustive and the Standard Integer Programming (SIP) approaches for the majority of the datasets used in the experiments. For the RASD approach, it is compared to other classification methods and it shows very competitive results in term of classification accuracy and reducts size.

As a conclusion, the proposed approaches have shown competitive and even better accuracy in most tested domains. The experiment results indicate that the proposed approaches as Rough classifiers give good performance across different classification problems and they can be promising methods in solving classification problems. Moreover, the experiments proved that the incremental vertical decomposition framework is an appealing method for knowledge reduction over large datasets within the framework of Rough Set based classification.



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

**PENDEKATAN SET ATRIBUT BERPEMBERAT DAN PEMECAHAN
BAGI PENGIRAAN PENGURANG**

Oleh

QASEM AHMAD AL-RADAIDEH

Julai 2005

Pengerusi : Profesor Madya Md. Nasir Sulaiman, Ph.D.

Fakulti: Sains Komputer dan Teknologi Maklumat

Penyelidikan ini adalah mengenai pengurangan pengetahuan berasaskan Set Kasar dan pengklasifikasian data dalam kerangka kerja perlombongan data. Bagi memudahkan pengklasifikasian berdasarkan Set Kasar, dua model utama bagi pengurang pengetahuan telah dicadang. Model pertama yang dicadangkan adalah pendekatan anggaran dalam pengiraan objek pengurang yang digunakan khusus untuk tujuan pengklasifikasian data. Pendekatan ini menekankan kepada penggunaan pemberat kepada setiap atribut di dalam set atribut. Pemberat-pemberat ini memberi petunjuk kepada kepentingan sesuatu atribut yang bakal dipertimbangkan di dalam pengurang. Pendekatan ini dinamakan POAB iaitu Pengurang Objek dengan Atribut Berpemberat. Satu variasi kepada pendekatan ini turut dicadangkan bagi mengira pengurang penuh. Variasi ini dinamakan sebagai PPAB bermaksud Pengurang Penuh dengan Atribut Berpemberat.

Model kedua yang dicadangkan melibatkan set data yang besar terutamanya dengan kuantiti atribut yang besar. Pendekatan ini menggunakan prinsip pemecahan set atribut secara berperingkat untuk menjana anggaran pengurang yang mewakili keseluruhan set data. Pendekatan yang dicadangkan ini dinamakan PPST bermaksud Pengurang dengan Pemecahan Set Atribut.

Pendekatan-pendekatan pengiraan pengurang yang dicadangkan dieksperimen dan dinilai secara menyeluruh. Proses penilaian adalah dalam dua aras: pertama adalah penilaian ke atas pendekatan yang dicadangkan berdasarkan Set Kasar di mana ketepatan pengklasifikasian digunakan sebagai ukuran penilaian. Kaedah penilaian bersilang 10-aras yang terkenal juga digunakan bagi menganggar ketepatan pengklasifikasian. Aras kedua penilaian digunakan untuk menilai pendekatan yang dikenali sebagai kaedah pengurang pengetahuan di mana saiz pengurang digunakan sebagai ukuran pengurangan.

Pendekatan-pendekatan ini dibandingkan dengan kaedah pengiraan pengurang yang lain dan termasuk lain-lain kaedah yang tidak berasaskan Set Kasar. Di dalam eksperimen, kami menggunakan pendekatan yang dicadangkan ke atas beberapa set data domain piawai daripada simpanan UCI. Keputusan eksperimen menunjukkan pencapaian yang sangat baik oleh pendekatan yang dicadangkan dalam proses pengklasifikasian dan pengurangan pengetahuan. Ketepatan pendekatan POAB melebihi pendekatan *Johnson* dalam kesemua set data. Ia juga menghasilkan ketepatan yang lebih baik jika dibandingkan dengan pendekatan *Exhaustive* dan *SIP*

dalam majoriti set data yang digunakan di dalam eksperimen. Bagi pendekatan PPSA, ianya juga telah dibandingkan dengan kaedah pengklasifikasian yang lain dan telah menunjukkan hasil keputusan yang kompetitif dari segi ketepatan pengklasifikasian dan saiz pengurang yang dijana.

Kesimpulannya, pendekatan-pendekatan yang dicadangkan telah menunjukkan ketepatan yang kompetitif, malah lebih baik apabila diuji menggunakan domain-domain ujian yang utama. Keputusan eksperimen menunjukkan pendekatan pengklasifikasi kasar yang dicadangkan berupaya memberi pencapaian yang baik dan menjanjikan hasil ke atas masalah-masalah pengklasifikasian. Tambahan pula, eksperimen telah membuktikan bahawa kerangka pemecahan menegak secara berperingkat adalah satu pendekatan yang menarik bagi pengurangan pengetahuan sekiranya menggunakan set data yang besar, dan ianya bernilai untuk digunakan di dalam kerangka pengklasifikasian berasaskan Set Kasar.

ACKNOWLEDGEMENTS

My thanks go firstly, as they should always be, to Allah, who blessed me with the ability to undertake and finally complete this work.

This work could not have been carried out without both direct and indirect help and support from my supervisor, Associate Professor Dr. Md Nasir Sulaiman. Thank you, Dr. Nasir, for introducing me to the topic of Rough Set theory, for support, and for opening doors.

Great thanks go to my supervisory committee members, Associate Professor Hj. Mohd Hasan Selamat and Associate Professor Dr. Hamidah Ibrahim for their valuable comments and fruitful discussions.

This research is partially supported by an IRPA fund. Thanks to University Putra Malaysia and the Malaysian government for the support.

My great thanks go to Yarmouk University, Jordan, for approving the leave to pursue my PhD study. Thanks to all colleagues there for their encouragements.

I will not forget to thank Prof. Dr. T.Y. Lin of San Jose State University, California, USA for his moral support and sharing with me some arguments and perspectives related to Data Mining.

My wife and my lovely daughters, *JOUD* and *DENA*, deserve my unending gratitude for their patience and understanding that make the duration of the study easier for me.

My Mother, Brother, and Sisters have always been a continuous source of support for which I am immensely grateful. They have been patient, encouraging and understanding. My thanks go to my brothers in law and uncles for their encouragements.

I also owe a great debt to my friends, mainly, Mohammad Abu Rahma and Saeed AlHazmi, for their support, for sharing me some enjoyable breaks from my work, and for many interesting conversations.

To everybody else, friends, colleagues, and relatives, thanks for their encouragements, cooperation, and moral support.

Qasem Al-Radaideh

July 2005



I certify that an Examination Committee met on 4th July 2005 to conduct the final examination of Qasem Ahmad Qasem Al-Radaideh on his Doctor of Philosophy thesis entitled "Attribute Set Weighting and Decomposition Approaches for Reduct Computation" in accordance with Universiti Pertanian Malaysia (Higher Degree) Act 1980 and Universiti Pertanian Malaysia (Higher Degree) Regulations 1981. The Committee recommends that the candidate be awarded the relevant degree. Members of the Examination Committee are as follows:

Mohamed Othman, PhD

Associate Professor
Faculty of Computer Science and Information Technology
Universiti Putra Malaysia
(Chairman)

Ramlan Mahmud, PhD

Associate Professor
Faculty of Computer Science and Information Technology
Universiti Putra Malaysia
(Internal Examiner)

Ali Mamat, PhD

Associate Professor
Faculty of Computer Science and Information Technology
Universiti Putra Malaysia
(Internal Examiner)

Razak Hamdan, PhD

Professor
Faculty of Information Science and Technology
Universiti Kebangsaan Malaysia
(External Examiner)



GULAM RUSUL RAHMAT ALI, PhD
Professor/Deputy Dean
School of Graduate Studies
Universiti Putra Malaysia

Date: 22 AUG 2005



I certify that an Examination Committee met on 4th July 2005 to conduct the final examination of Qasem Ahmad Qasem Al-Radaideh on his Doctor of Philosophy thesis entitled "Attribute Set Weighting and Decomposition Approaches for Reduct Computation" in accordance with Universiti Pertanian Malaysia (Higher Degree) Act 1980 and Universiti Pertanian Malaysia (Higher Degree) Regulations 1981. The Committee recommends that the candidate be awarded the relevant degree. Members of the Examination Committee are as follows:

Mohamed Othman, PhD

Associate Professor
Faculty of Computer Science and Information Technology
Universiti Putra Malaysia
(Chairman)

Ramlan Mahmud, PhD

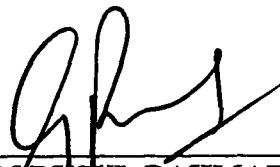
Associate Professor
Faculty of Computer Science and Information Technology
Universiti Putra Malaysia
(Internal Examiner)

Ali Mamat, PhD

Associate Professor
Faculty of Computer Science and Information Technology
Universiti Putra Malaysia
(Internal Examiner)

Razak Hamdan, PhD

Professor
Faculty of Information Science and Technology
Universiti Kebangsaan Malaysia
(External Examiner)



GULAM/RUSUL RAHMAT ALI, PhD
Professor/Deputy Dean
School of Graduate Studies
Universiti Putra Malaysia

Date: 22 AUG 2005



This thesis was submitted to the Senate of Universiti Putra Malaysia and has been accepted as fulfilment of the requirements for the degree Doctor of Philosophy. The members of the Supervisory Committee are as follows:

Md Nasir Sulaiman, PhD
Associate Professor
Faculty of Computer Science and Information Technology
Universiti Putra Malaysia
(Chairman)

Mohd Hasan Selamat
Associate Professor
Faculty of Computer Science and Information Technology
Universiti Putra Malaysia
(Member)

Hamidah Ibrahim, PhD
Associate Professor
Faculty of Computer Science and Information Technology
Universiti Putra Malaysia
(Member)



AINI IDERIS, Ph.D.
Professor / Dean
School of Graduate Studies
Universiti Putra Malaysia

Date: 08 SEP 2005



DECLARATION

I hereby declare that the thesis is based on my original work except for quotations and citations which have been duly acknowledged. I also declare that it has not been previously or concurrently submitted for any other degree at UPM or other institutions.



QASEM A. Q. AL-RADAIDEH

Date: 20/8/2005

TABLE OF CONTENTS

	Page
DEDICATION	ii
ABSTRACT	iii
ABSTRAK	v
ACKNOWLEDGEMENTS	viii
APPROVAL	ix
DECLARATION	xi
LIST OF TABLES	xv
LIST OF FIGURES	xviii
LIST OF ABBREVIATIONS AND NOTATIONS	xix
CHAPTER	
1 INTRODUCTION	
1.1 Background	1
1.2 Problem Statement	3
1.3 Objective of the Research	5
1.4 Scope of the Research	6
1.5 Research Methodology	6
1.6 Contributions of the Research	8
1.7 Organization of the Thesis	9
2 LITERATURE REVIEW	
2.1 Introduction	12
2.2 Data Mining and Knowledge Discovery	13
2.2.1 Knowledge Discovery Phases	14
2.2.2 Data Mining Tasks	17
2.3 Data Preprocessing	19
2.4 Classification in Data Mining Framework	20
2.4.1 Induction and Classification Methods	21
2.4.2 Simplicity First Principle and the OneR Classifier	22
2.4.3 Decision Tree based Classification Methods	23
2.4.4 Classification Evaluation	24
2.4.5 Methods for Accuracy Estimation	25
2.4.6 Confusion Matrix	28
2.4.7 Rules and Decision Making	30
2.5 Data Mining and Large Data	30
2.5.1 Data Dimension Reduction	31
2.5.2 Decomposition in Data Mining	33
2.6 Rough Set Theory : Preliminaries and Extensions	33
2.6.1 Rough Set Theory Preliminaries	34
2.6.2 Information and Decision Systems	34
2.6.3 Indiscernibility Relation	36
2.6.4 Reduct and Core	37
2.7 Rough Set Theory Extensions	39



2.7.1	Discernibility Concepts	40
2.7.2	Focused Decision Table and Focused Reduct	44
2.7.3	Attribute Set Ranking	45
2.8	Reduct Computation	45
2.8.1	The Problem and Solutions	45
2.8.2	Importance of Reduct	47
2.9	Review of Reduct Computation Related Work	49
2.9.1	Discernibility Matrix based Heuristic Approaches	50
2.9.2	Johnson Approximation for Reduct Computation	56
2.10	Rule Set Generation	58
2.11	Rough Set based Classification	60
2.12	Data Decomposition for Data Mining	65
2.13	Datasets Reduction by Decomposition Principle	66
2.13.1	Attribute Set Decomposition	66
2.13.2	Object Set Decomposition	68
2.14	Decomposition Structure	69
2.15	Review of Some Attribute Set Decomposition Approaches	69
2.15.1	Decomposition for Dedicated Applications	70
2.15.2	Decomposition for Decision Making	72
2.15.3	Decomposition for data Classification	73
2.15.4	Decomposition for Reduct Computation	74
2.16	Summary	75

3 ATTRIBUTE SET WEIGHTING FOR REDUCT COMPUTATION

3.1	Introduction	78
3.2	The Proposed Approximate Reducts Computation Model	79
3.2.1	Attribute Set Weighting	79
3.2.2	Attribute Local Weighting	81
3.2.3	Attribute Global Weighting	82
3.2.4	Attribute Cardinality	82
3.2.5	An Overview of the Algorithm	83
3.3	Illustrative Example for Object Reduct Computation	86
3.4	Improved Discernibility Matrix Structure	90
3.5	Full Reduct Computation	92
3.6	Illustrative Example for Full Reduct Computation	93
3.7	The Complexity and Implementation of the Approach	94
3.8	Building the Rough Classifier	95
3.9	Summary	100

4 ATTRIBUTE SET DECOMPOSITION FOR REDUCT COMPUTATION

4.1	Introduction	102
4.2	Notations and Definitions	103
4.3	Reduct by Incremental Attribute Set Decomposition	104
4.3.1	Overview of the Proposed Approach	105
4.3.2	Simple Illustrative Example	109
4.4	The Extreme Case Decomposition	111
4.5	Simple Illustrative Example for the Extreme Case	112
4.6	Summary	114



5	RESULTS AND DISCUSSIONS	
5.1	Introduction	115
5.2	Experiments Setup	116
5.3	Datasets Used in Experiments	119
5.4	Dataset Preprocessing	119
5.5	Experiments for the Object Reduct Computation Approach (ORAW)	121
5.5.1	Comparison with Johnson Approach	122
5.5.2	Comparison with other Reduct Computation Approaches	125
5.5.3	Comparison with other Classification Methods	127
5.6	Experiments for the Full Reduct Computation Approach (FRAW)	130
5.7	Evaluating the Incremental Decomposition Approach (RASD)	131
5.7.1	Comparison of RASD with other Classification Methods	137
5.7.2	Experiments over Large Datasets	138
5.7.3	Comparison with other Decomposition based Methods	141
5.8	Summary	142
6	CONCLUSIONS AND FUTURE WORK	
6.1	Introduction	144
6.2	Concluding Remarks	145
6.3	Future Work and Extensions	148
	BIBLIOGRAPHY	150
	APPENDIX A	159
	BIODATA OF THE AUTHOR	173
	LIST OF PUBLICATIONS RELATED TO THIS THESIS	174



LIST OF TABLES

Table		Page
2.1a	An Example of a Decision System S	36
2.1b	The Decision System S in Numeric Form	36
2.2	The Discernibility Matrix of S	41
2.3	The Discernibility Matrix Modulo of S	42
2.4	Object Discernibility Functions of S	43
3.1	An Example of a Decision System S	87
3.2	The Discernibility Matrix (DM) and the Reduct Set of S	87
3.3	DM after Removing Entries that Contain a_1	88
3.4	DM after Removing Entries that Contain a_5	89
3.5	The Final DM and Reduct Set after Ending the Algorithm	90
3.6	The Discernibility List Structure of the DMM for S	91
3.7	Lower Triangle of DMM of the Decision System S	93
3.8	The DMM after Removing Entries that Contain a_4	94
3.9	The Set of Rules Generated for Decision System S	99
4.1	The List Structure Contents of the Decision System SI	111
4.2	The List Structure Contents of the Decision System S	113
5.1	Characteristics of the Datasets	120
5.2	The Dataset Preprocessing	121
5.3	Datasets Classification Accuracy using the ORAW Approach	122
5.4	Datasets Classification Accuracy using the Johnson Reducer	123
5.5	Comparison of the Accuracy of ORAW and Exhaustive Approaches	126
5.6	Comparison of the Accuracy of ORAW and SIP Approaches	127
5.7	Comparison of ORAW with the Popular Classification Approaches	128
5.8	ORAW t -test Results Comparing to other Classification Approaches	129



Table	Page	
5.9	Reduct Computation using FRAW Approach	130
5.10	Datasets Accuracy Obtained by RSAD	132
5.11	Comparison of the Accuracy of RASD with FSR Approach	134
5.12	Comparison of the RASD and Zhang Approaches	135
5.13	Comparison of Different Reduct Computation Methods for the Zoo Dataset	136
5.14	Comparison of RASD with the Popular Classification Approaches	137
5.15	Experimenting RASD Approach over the DNA dataset	138
5.16	Comparing RASD Approach with Other Classification Methods over the DNA Dataset	139
5.17	Comparison of RASD and DIFN	141
5.18	Summary Comparison of ORAW and other Classification Approaches	143
A.1	Sample of the Australian Credit Card Dataset after Preprocessing	159
A.2	Sample of the Chess Dataset after Preprocessing	160
A.3	Sample of the Cleveland Dataset after Preprocessing	161
A.4	Sample of the Glass Dataset after Preprocessing	161
A.5	Sample of the Heart Disease Dataset after Preprocessing	162
A.6	Sample of the Hepatitis Disease Dataset after Preprocessing	162
A.7	Sample of the Iris Dataset after Preprocessing	163
A.8	Sample of the Letter Dataset after Preprocessing	163
A.9	Sample of the Lymphography Dataset after Preprocessing	164
A.10.a	Sample of the Lung Cancer Dataset after Preprocessing	164
A.10.b	(Continue) Sample of the Lung Cancer Dataset after Preprocessing	165
A.11.a	Sample of the Sonar Dataset after Preprocessing	166
A.11.b	(Continue) Sample of the Sonar Dataset after Preprocessing	167



Table		Page
A.12	Sample of the Soya Beans (large) Dataset after Preprocessing	168
A.13	Sample of the Vehicle Dataset after Preprocessing	169
A.14	Sample of the Zoo Dataset after Preprocessing	169
A.15.a	Sample of the DNA Dataset after Preprocessing	171
A.15.b	(Continue) Sample of the DNA Dataset after Preprocessing	171
A.15.c	(Continue) Sample of the DNA Dataset after Preprocessing	172
A.15.d	(Continue) Sample of the DNA Dataset after Preprocessing	172



LIST OF FIGURES

Figure		Page
2.1	The Knowledge Discovery Process	15
2.2	Data Classification Process	21
2.3	The Confusion Matrix Structure	29
2.4	Johnson Approximation Algorithm	57
2.5	The Pseudo-code of Decision Rules Generation Algorithm	59
2.6	Rough Classification Framework	62
3.1	The Pseudo-code of the Object Reduct Computation Algorithm (ORAW)	84
3.2	Building the List Structure of the Discernibility Matrix Modulo	91
3.3	The Pseudo-code of the Full Reduct Computation Algorithm (FRAW)	92
3.4	Flow Diagram of the Rough Set Based Classification Process	96
3.5	Top Level Pseudo-code of the Rule Extraction Phase	96
4.1	The Flow Diagram of the Reduct Computation by Incremental Decomposition Approach (RASD)	106
4.2	The Pseudo-code of the RASD Algorithm	106
4.3	The Overlapping of Partitions	107
4.4	An Example of the Incremental Decomposition Approach ($K=4$)	110
4.5	An Example of the Extreme Case Decomposition ($K=n$)	112
5.1	Chart Comparison of ORAW and Johnson Approaches	123
5.2	The <i>t-test</i> Results for Means of the Datasets	129
5.3	Chart Comparison of ORAW and other Classification Approaches	129



LIST OF ABBREVIATIONS AND NOTATIONS

<i>A</i>	Set of conditional attributes
ACC	Classification ACCuracy
AVCV	Attribute Value Cardinality Vector
c_{ij}	Entry of Discernibility Matrix
CV	Cross Validation
<i>d</i>	The decision attribute in a decision system
DL	Discernibility List
DM	Discernibility Matrix
DMM	Discernibility Matrix Modulo
DS	Decision System
FRAW	Full Reduct by Attribute set Weighting
GAWV	Global Attribute Weight Vector
<i>gw</i>	global weight
IND	INDiscernibility relation
IS	Information System
KDD	Knowledge Discovery in Database
LAWV	Local Attribute Wight Vector
<i>lw</i>	local weight
OneR	One Rules
ORAW	Object Reducts by Attribute set Weighting
RASD	Reduct by Attribute Set Decomposition
RED	REDuct Set
ROSETTA	ROugh SET Toolkit for Analysis of data
RSES	Rough Set Exploration System
RSESlib	Rough Set Exploration System library
<i>U</i>	Universe of objects
<i>v</i>	An attribute value
<i>V_a</i>	The value set of an attribute
<i>vcw</i>	value cardinality weight



CHAPTER 1

INTRODUCTION

1.1 Background

Due to the explosion of data in our modern society, most organizations have large databases that contain a wealth of undiscovered, yet valuable information. To gain benefits from the collected information and to discover the valuable knowledge, it needs to be analyzed. This leads to a need for methods and ways to aid or substitute humans in the process of knowledge discovery from large datasets. Knowledge discovery and data mining methodologies have been introduced as methods for bridging the knowledge gap between information gathered and information analyzed (Han & Kamber, 2001; Cios *et al.*, 1998; Fayyad *et al.*, 1996b, 1996c). Analogous to the mining in the real world, data mining is that, with the computer, we can automatically find the “information gold nuggets” or “diamonds” by sifting out enormous quantities of data-debris from our database.

Data mining is a promising and an interdisciplinary research area spanning several disciplines such as database, machine learning, artificial intelligence, intelligent information systems, statistics, data warehousing and knowledge acquisition in expert systems. Data mining has evolved into an important and active area of research because of theoretical challenges and practical application associated with the problem of discovering interested or previously unknown knowledge from very



large real-world databases. With data mining we can simply let data “speak for itself”.

There are several tasks in data mining and the most common in the literature is classification, which is a form of data analysis that can be used to extract models describing important data classes. The classification task concentrates on predicting the value of the decision class for an object among a predefined set of classes’ values given the values of some given attributes for the object.

In the literature many classification approaches have been proposed and implemented by researchers, such as, decision tree based classification, statistical classification, neural network based classification, genetic algorithms classifiers and Rough Set based classification (Cios *et al.*, 1998; Bazan *et al.*, 2000). Classification has a wide range of applications, including scientific experiments, medical diagnosis, credit approval, etc.

Rough Set theory is a mathematical tool developed as a formal method to turn data into knowledge (Pawlak, 1991). The two main applications of the classical Rough Sets theory are in attribute reduction and classification. Rough Set based classification is inspired by the concepts of the Rough Set theory with a primary goal to extract rules from data represented in a decision system. According to Pawlak (1991), the notion of classification is central to the theory.

A very important issue in data mining is the data redundancy where not all knowledge presented to the data mining task in an information system is necessary to

describe it (Pawlak, 1991; Kohavi & Frasca, 1994; Zhang & Yao, 2004; Zhang *et al.*, 2003; Lin & Yin, 2004; Hu *et al.*, 2000; Boussouf & Quafafou, 2001). It is often the case where some of attributes or some of attributes values are superfluous. Rough Sets theory provides the *reduct* concept for data reduction as preprocessing step of data analysis. A reduct is defined as the minimal attribute set preserving classification power of the original information system with the full set of attributes.

The reduct concept of the theory is a fundamental concept towards rule extraction. The concept enables us to discard functionally the redundant information and guarantees that the attributes that do not contribute to the classification are removed. The process of finding reduct is a fundamental step in applying the Rough Set theory for the data classification task. Based on the reduct concept, the rules generated by the classifiers are expected to be more concise than if generated over the original dataset (Pawlak, 1998; Komorowski *et al.*, 1999).

1.2 Problem Statement

Data classification problem is a well known problem in the area of knowledge discovery. In applying Rough Set theory as a classification framework, the problem of computing reducts as a knowledge reduction method, is without doubt the most complex and computer-intensive step in Rough Set data analysis (Pawlak, 1998). The problem of computing all reducts is known to belong to a theoretical class of problems that, informally, requires an amount of computation that grows exponentially with the size of the problem. The problem size is dominated by the number of attributes and objects involved.

Several approximation and heuristic methods have been proposed but there are no universal solutions and no accredited best heuristic method (Kuo & Yajima, 2003). According to Lin & Yin (2004), Kuo & Yajima (2003), and Wang & Chen (2004), so far, the problem of reduct computation stills an open research area in Rough Set theory particularly for large datasets with large number of attributes.

Generally, most of the available heuristic approaches use the discernibility matrix concept and a weighting mechanism to evaluate the significance of the attributes to be considered in the reduct Zhang *et al.* (2003). The available weighting mechanisms may lead to consider some attributes with less importance which eventually lead to low classification accuracy. In addition, some of the available approaches have limitations in handling large amount of datasets particularly with large number of attributes (Bakar *et al.*, 2002; Zhengren *et al.*, 2004).

In the available approaches, the most used weight for attributes is the number of occurrences in the discernibility matrix and when several attributes have the same weight a random choice is used. This may allow less significant attributes to be a member of the reduct which lead to low classification accuracy.

Johnson reducer (Nguyen & Nguyen, 1996, Ohrn, 1998) uses the attribute frequency in the discernibility matrix to measure the significance of attributes to be considers in the reduct. A random choice is made when several attributes have the same significance. Hu *et al.* (2000) use the attribute frequency and entry length in discernibility matrix as measures for the significance of attributes. The same

