



AIN SHAMS UNIVERSITY
FACULTY OF COMPUTER & INFORMATION SCIENCES
Department of Information Systems

HIGH PERFORMANCE DATA MINING IN DISTRIBUTED DATABASES

**Thesis submitted as a partial fulfillment of the requirements for the
degree of Master of Science**

In

Computer and Information Sciences

By

Mahmoud Fouad Anwar Darwish

B.Sc. in Computer and Information Sciences, Information Systems
Department, Ain Shams University

Under Supervision of

Prof. Dr Hossam El Deen Faheem

Professor
Computer Systems Department,
Faculty of Computers and Information Sciences,
Ain Shams University, Egypt

Prof. Dr Nagwa Badr

Professor
Information Systems Department,
Faculty of Computers and Information Sciences,
Ain Shams University, Egypt

Prof. Rania El Gohary

Assistant Professor
Information Systems Department,
Faculty of Computers and Information Sciences,
Ain Shams University, Egypt

Abstract

The massive volume of data generated on daily basis decreases the ability of current data mining techniques to generate knowledge in a short time. The constant change in data requires constant updating of the existing patterns. It is computationally intensive to repeat the knowledge discovery process on the whole databases with every update. Therefore, there is a need to enhance the performance of association rules mining methodologies when dealing with incremental updates.

In order to enhance the performance of incremental association rules mining, this thesis focus on the utilization of current hardware and software advances in high-performance computing. This thesis proposes a distributed incremental association rules mining approach based on MPI. In addition, the thesis also proposes a hybrid incremental mining approach based on OpenMP and MPI to work in high performance computing environments. In order to reduce the need to reprocess the entire database, this thesis depends on pre-large and negative borders approaches.

To evaluate the applied approaches, this thesis considered the output accuracy, processing time and the acceleration as our primary evaluation metrics. In fact, experimental results have proved that our distributed method reduces processing time by 40% when compared to

serial existing approach and our hybrid approach reduces processing time by 19% when compared to distributed approach.

Acknowledgments

I would like to Thank God for his blessings and support to finalize this thesis.

I would like to thank Mrs. Fahima Moustafa, my mother, the real hero behind this thesis and any success in my life. I am so lucky to have a wonderful mother like you.

I would like to thank my supervisors for their support and help to accomplish my master. I wouldn't be able to do this without your guidance and support. A Special Thank you goes to Dr. Rania El Gohary, Dr. Nagwa Bader, and Prof. Hossam El-Deen Faheem.

I also thank the teaching assistants at Faculty of Computers and Information Science for their support during performing experiments at high performance computing laboratory. Thank you. I would also like to Thank Mohamed Gawish for his continuous support throughout this thesis.

I would like to thank Ain Shams University, Faculty of Computer and Information Science for providing the appropriate atmosphere and facilities during my master. I always feel so proud for being a graduate from the faculty of Computer & Information Science, Ain Shams University.

Table of Contents

CHAPTER	PAGE
List of Tables	vii
List of Figures.....	viii
List of Abbreviations	ix
Chapter 1: Introduction	1
1.1 Overview	1
1.2 Problem Definition.....	4
1.3 Motivation	4
1.4 Thesis Objectives	4
1.5 Thesis Outlines.....	5
Chapter 2: Background and Literature Review	6
2.1 Association Rules Mining.....	6
2.1.1 Association Rules Mining Problem.....	6
2.1.2 Apriori.....	8
2.1.3 Count Distribution	9
2.1.4 Data Distribution.....	9
2.1.5 Candidate Distribution	10
2.1.6 Fast Distributed Mining	10
2.1.7 Distributed Decision Miner	11
2.1.8 Sampling.....	11
2.1.9 Distributed-Sampling.....	12
2.1.10 Optimized distributed association rule mining algorithm.....	12
2.1.11 Secure Incremental Maintenance of Association Rules.....	13
2.1.12 Frequent Pattern-Growth	13
2.1.13 Synthesizing global patterns from local patterns	14
2.1.14 Discussion of Association Rules Mining Methodologies.....	16
2.2 Incremental Association Rules Mining.....	18
2.2.1 Fast Update	19
2.2.2 Fast Updated 2	20
2.2.3 Incremental Frequent Pattern - Growth.....	21
2.2.4 Maintaining Association Rules with Apriori Property.....	21
2.2.5 Negative Border	22
2.2.6 Pre-Large Approach.....	23
2.3 High Performance Computing	26
2.3.1 Message Passing Interface	26
2.3.2 OpenMP.....	28
2.3.3 Hybrid MPI/OpenMP	29
Chapter 3: Proposed Approach.....	31
3.1 Distributed Approach of Incremental Association Rules Mining.....	31
3.1.1 Overview.....	31

3.1.2 Proposed Approach High Level Design	32
3.1.3 System Components and Implementation	35
3.1.4 Proposed Approach Detailed Design	38
3.1.5 Illustrated Example on the proposed approach	43
3.2 Hybrid approach of Incremental ARM using MPI/OPENMP	61
3.2.1 Proposed Approach High Level Design	62
3.2.2 Proposed Approach Detailed Design	64
Chapter 4: Experimental Results	67
4.1 Distributed Approach Experimental Results	67
4.2 Hybrid Approach Experimental Results	71
Chapter 5: Conclusion and Future Work.....	79
5.1 Conclusion	79
5.2 Future Work.....	80
References	81

List of Tables

Table	Page
Table 2.1 Association Rules Mining Techniques Comparison.....	17
Table 3.1 CairoDB Original Database	43
Table 3.2 KafreElShiekhDB Original Database	43
Table 3.3 Original Large Itemsets	44
Table 3.4 Original pre-large Itemset	45
Table 3.5 Original Small Itemsets	45
Table 3.6 new added incremental records.....	45
Table 3.7 Consolidated Original Database Mapping.....	48
Table 3.8 Consolidated New Added Records Database Mapping.....	48
Table 3.9 Consolidated Simplified View of New Added Records Database Mapping	49
Table 3.10 First Incremental Horizontal Partition	49
Table 3.11 Second Incremental Horizontal Partition.....	49
Table 3.12 1st itemsets candidate list - worker[1] result	50
Table 3.13 1st itemsets candidate list - worker [2] result	50
Table 3.14 1st global itemsets candidate list - master combination result.....	51
Table 3.15 Updated Support Count for 1st itemsets candidates - Master.....	52
Table 3.16 Updated 1 st candidates Large/Pre-Large/Small lists - Master process	53
Table 3.17 1st itemsets candidates Original vs New Large/Pre-Large/Small lists	53
Table 3.18 First itemsets candidate list Input for Worker Process [1].....	54
Table 3.19 First itemsets candidate list Input for Worker Process [2].....	54
Table 3.20 2nd itemsets candidate list - worker [1] result.....	55
Table 3.21 2nd itemsets candidate list - worker [2] result.....	55
Table 3.22 2nd global itemsets candidate list - master combination result	55
Table 3.23 Updated Support Count for 2nd itemsets candidates – Master.....	56
Table 3.24 Updated 2nd Itemsets candidates Large/Pre-Large/Small lists	57
Table 3.25 2nd itemsets candidates' original vs new large/pre-large/small lists.....	57
Table 3.26 2nd itemsets candidate list Input for worker processes	58
Table 3.27 3rd itemsets candidate list - worker [1] result.....	59
Table 3.28 3rd itemsets candidate list - worker [2] result.....	59
Table 3.29 Updated Support Count for 3rd itemsets candidates – Master	60
Table 3.30 Updated 3rd Itemsets candidates - Master process	60
Table 3.31 3rd itemsets candidates' original vs new incremental small list.....	60

List of Figures

Figure	Page
Figure 1.1 Today's dataset characteristics	1
Figure 2.1 Frequent Itemset Example	7
Figure 2.2 Animesh approach for generating global patterns from local patterns.....	15
Figure 2.3 Incremental Mining Process	19
Figure 2.4 FUP Cases	20
Figure 2.5 Association rules maintenance cases arising in records modifications	26
Figure 2.6 MPI Distributed Memory Architecture	27
Figure 2.7 MPI Distributed & Shared Memory Architecture	28
Figure 2.8 OpenMP Uniform Memory Access Architecture.....	29
Figure 2.9 OpenMP Non Uniform Memory Access	29
Figure 2.10 Hybrid Architecture MPI/OpenMP	30
Figure 3.1 High Level Diagram for Distributed Incremental Mining Approach.....	32
Figure 3.2 Data Flow Parallel & Distributed Incremental ARM approach	34
Figure 3.3 System Components of Proposed Distributed Incremental ARM.....	36
Figure 3.4 Proposed Distributed System Flowchart	39
Figure 3.5 Slave/Worker Process Pseudocode.....	40
Figure 3.6 Master/Manager Pseudocode.....	42
Figure 3.7 Worker Node Pseudocode	63
Figure 3.8 Manager- Worker Communication Diagram.....	64
Figure 4.1. Distributed Approach Execution time comparison 1	69
Figure 4.2. Distributed Approach Execution time comparison 2	70
Figure 4.3. Execution Time Comparison on Linux Cluster – Orders Database	73
Figure 4.4. Execution Time Comparison on Windows Cluster - Orders Database	74
Figure 4.5. Execution Time Comparison on Linux Cluster – Telecom IVR Database	75
Figure 4.6. Execution Time Comparison on Windows Cluster, Telecom IVR Database	75
Figure 4.7. Hybrid vs Distributed Approach Speedup on Windows Cluster.....	76
Figure 4.8. Hybrid vs Distributed Approach Speedup on Linux Cluster.....	77

List of Abbreviations

ARM	Association Rules Mining
DDM	Distributed Decision Miner
Di	Database number i
FDM	Fast Distributed Mining
FP-Growth	Frequent Pattern Growth
FUP	Fast Update
FUP2	Fast Updated2
HPC	High Performance Computing
ICSC	Incremental candidate support count
IDC	New incremental database count
LP	Local Pattern
LS	Large support
MAAP	Maintaining Association Rules with Apriori Property
MPI	Message Passing Interface
NLSC	New large minimum support count
NPSC	New pre-large minimum support count
OCSC	Original candidate support count
ODAM	Optimized Distributed Association Rule Mining
ODC	Original database transactions count
OpenMP	Open Multi-Processing
PLS	Pre large Support

SIMDAR	Secure Incremental Maintenance of Distributed Association Rules
SLP	Suggested Local Pattern
UCSC	updated candidate support count in both original and incremental databases

Chapter 1: Introduction

1.1 Overview

In our present time, data volume is increasing dramatically, sources of data are becoming decentralized, and the rate of data updates is tremendously speeding up. Today's data characteristics make it hard for humans to understand the big picture or get accurate insights about these massive amounts of data. Figure 1.1 shows the main challenging characteristics of today's dataset. The first challenge is the massive volume of data generated every second by various systems. The second one is velocity as the speed of data generation is increasing every second. The third challenge is variety of datatypes as nowadays there are different types of structured, unstructured and semi-structured data sources. These problems makes it necessary to have an automated analytical method that could deal with today's data and generate useful insights in a short time.

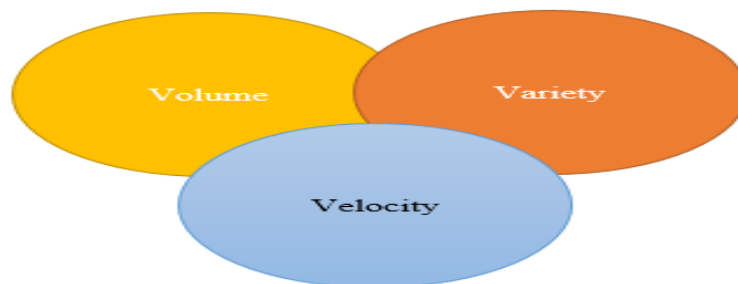


Figure 1.1 Today's dataset characteristics

Data mining provides the ability to discover hidden knowledge from data. It comprises different techniques which reveal the useful patterns that are implicit in the data [1-2]. Through data mining enable us to extract useful patterns and relationships between data using the association rules mining method. It also enables us to group similar data together using the clustering technique. In addition, data mining makes it possible for us to predict the future through analyzing current data [3], [4]. With data mining one can also analyze data in less time especially when compared to the time consumed to examine it manually. Consequently, data mining is revolutionary as it provides hope to develop automated tools to process data to generate useful knowledge [5].

Association rules mining is considered an important method of data mining and a well-studied problem as well. It is used to predict the occurrence of an item depending on the occurrence of another item [6], [7]. The problem originates in supermarket sales systems. For example, the rule {bread, cheese} \rightarrow {eggs} predicts that a customer is going to buy eggs in case he buys bread and cheese. Such rules are very important in empowering decision makers to take the right decision in promotional pricing and marketing. Besides market basket analysis, Association rules mining is used in different areas like detecting intrusion [8-10]; mining usage of the web [11-13]; medical diagnosis [14-16]; and others [17-19].

With the continual and constant increase of data, the generated association rules based on old data may become invalid and outdated. To maintain the existing rules, Association rules mining should be applied again on the whole dataset after every change. This process is called incremental association rules mining [20-24]. Researchers proposed different approaches for handling incremental association rules mining, however, performance is still a hot area of research. Some could handle incremental mining but require more than one database scan for original data in some cases like FUP and FUP2. Others methods like FP-Growth could take one scan but can't fit in memory. [25-30]

High performance computing is the use of parallel processing for running applications efficiently and quickly. The latest advances in technology are pushing towards parallel processing. Many computing cores could exist in single microprocessor. Current hardware clusters support hybrid memory nature with a distributed memory across all nodes and non-uniform memory access in each node. All cluster nodes are connected together through a high speed network. Various high performance programming modules that could be used in parallel processing has been introduced lately like MPI and OpenMP. Developing a hybrid data mining applications that could make use of cluster's hybrid memory nature is a promising field that could improve performance dramatically [31].

1.2 Problem Definition

The massive amount of data generated in the real-world applications has a significant effect on the speed of processing. Therefore, the speed of maintaining association rules decreases whenever data increases. There are multiple problems in improving association rules discovery and maintenance processes. The first problem is to reduce any need to reprocess the whole database whenever data changes. The second problem is to reduce the execution time for processing newly added dataset. The third problem is to allow incremental mining approach to process multiple data concurrently.

1.3 Motivation

Rapid development realms of parallel computing paradigms and parallel processing could have been more efficiently utilized to tackle the problem of incremental association rule mining in order to achieve much more advantageous operation.

1.4 Thesis Objectives

The objective of this thesis includes:

A) Providing a high performance distributed incremental association rules mining approach that can deal with today's dataset in an efficient and quick way.

B) Using high-performance computing environments to reduce processing time for incremental association rules mining.

C) Reducing the need to reprocess the whole database to maintain association rules.

1.5 Thesis Outlines

This thesis is structured as follows:

Chapter 1 is the Introduction to the thesis. Chapter 2 presents the background and literature review. Chapter 3 focuses on the proposed approach. Chapter 4 deals with experimental results and findings. Chapter 5 presents the Conclusion and future work.