



ASCF: Apriori Algorithm on Spark Based on Cuckoo Filter Structure

By

Bana Ahmad Alrahwan

A Thesis Submitted to the
Faculty of Engineering at Cairo University
In Partial Fulfillment of the
Requirements for the Degree of
MASTER OF SCIENCE
In
Computer Engineering

ASCF: Apriori Algorithm on Spark Based on Cuckoo Filter Structure

By

Bana Ahmad Alrahwan

A Thesis Submitted to the Faculty of Engineering at Cairo University In Partial Fulfillment of the Requirements for the Degree of MASTER OF SCIENCE In **Computer Engineering**

Under the Supervision of

Prof. Dr. Elsayed Hemayed

Dr. Mona F. Ahmed

Professor Computer Engineering

Assistant Professor Computer Engineering Faculty of Engineering, Cairo University

Faculty of Engineering, Cairo University

> FACULTY OF ENGINEERING, CAIRO UNIVERSITY GIZA, EGYPT 2018

ASCF: Apriori Algorithm on Spark Based on Cuckoo Filter Structure

By

Bana Ahmad Alrahwan

A Thesis Submitted to the
Faculty of Engineering at Cairo University
In Partial Fulfillment of the
Requirements for the Degree of
MASTER OF SCIENCE

In Computer Engineering

Approved by the Examining Committee	
Prof. Dr. Elsayed Hemayed	Thesis Main Advisor
Prof. Dr. Magda Fayek	Internal Examiner
Prof. Dr. Reda Abdelwahab (Faculty of Computer and Information, IT Department)	External Examiner

Engineer's Name: Bana Ahmad Alrahwan

Nationality: Syrian **Date of Birth:** 3/01/1978

E-mail: bana.alrahwan-1@eng1.cu.edu.eg

Phone: 01092799738

Address: 29 Fahmy Street, Almanial

Registration Date:1/03/2014Awarding Date:..../..../2018Degree:Master of ScienceDepartment:Computer Engineering



Supervisors:

Prof. Dr. Elsayed Hemayed

Dr. Mona F. Ahmed

Examiners: Prof. Dr. Elsayed Hemayed (Thesis Main Advisor)

Prof. Dr. Magda Fayek (Internal Examiner)
Prof. Dr. Reda Abdelwahab (External Examiner)

(Faculty of Computer and Information, IT Department)

Title of Thesis:

ASCF: Apriori Algorithm on Spark Based on Cuckoo Filter Structure

Key Words:

Data mining, Big Data, Apache Spark, Apriori algorithm, Cuckoo filter.

Summary:

The Apriori algorithm is one of the most basic techniques that are used to discover frequent patterns in dataset. Apriori is iterative and works sequentially. It generates candidate sets having all possible combinations for frequent itemsets that are generated from the previous iteration and comparing each combination of items with every transaction record in every iteration. Thus, Apriori algorithm is not efficient and gets computationally more expensive as the data size is increased. The rapid growth of data necessitates running the data intensive algorithms in parallel distributed environment to achieve convenient performance. Many approaches have been proposed to solve the Apriori major drawbacks that severely degrade the performance as the datasets get larger which is a common feature in Today's data. In this thesis, Apriori Algorithm on Spark based on Cuckoo Filter structure (ASCF) is introduced. ASCF succeeds in removing the candidate generation step from Apriori algorithm to reduce computational complexity and avoid costly comparisons. The proposed algorithm is implemented on spark in-memory processing distributed environment to reduce processing time. The ASCF offers great improvement in performance over other implementation approaches of Apriori algorithm based on spark.

Acknowledgments

Today, I have completed my master thesis, I would like to take this opportunity firstly and foremost thank **Allah** for supplied me with the courage, guides me to the right path and has helped me to complete this thesis. Without Him, I would not have had the wisdom or the physical ability to do so.

I cannot forget the ideal man of the world and most respectable personality for whom **Allah** created the whole universe, **Prophet Mohammed** (Peace Be Upon Him).

This thesis would not have been possible without the help, encouragement and support of a number of wonderful individuals during the journey of preparing this thesis — So, I have to give them all thanks and appreciation.

I would like to thank my supervisor **Prof. Dr. Elsayed Hemayed** for his continuous efforts and concern. He has been always willing to help and give his best suggestions.

Thanks Dr. Elsayed

I would like to thank my gratitude to my supervisor **Dr. Mona Farouk** for suggesting the research path, her consistent guidance, advices and optimism. She has given me a confidence in my abilities, and provided me a strong supports to successfully complete my thesis. She has been very keen on my progress in this thesis. I greatly appreciate everything she has done for me.

Thanks Dr. Mona

I would like to express my deep and sincere gratitude to my family for their continuous and unparalleled love, help and support they provided me through my entire life. I dedicate my achievement to them.

My Beloved father, your absence had broken my heart and lost me the meaning of life, but your spirit still forever dwell in my heart. Your love surrounds me with every breath I take. You had taught me bravery, patience and sacrifice. You had always been my best friend, my greatest teacher. Today, I complete my career which he encouraged me to started it. I miss him so much, and I wish if he was with me. My father will remain a crown to be proud of all my life. You were and will always be the greatest father.

My Dear mother, I can barely find the words to express all the wisdom, love and support you've given me. Your being in my life is the secret of my success and happiness. Thanks for everything you have given me. Thank you for being my mother. **I Love You.**

My Brothers, Engineer Mustafa who is like my father in his tenderness and his sacrifice and **Dr. Mohammad** who accompanied me during my studies journey, he was my best friend. Thank you for your continuous support.

My Sisters, Fatimah, Engineer Lina, Razan and Dr. Ghofran, thank you for your encouraging me towards success.

Thanks **My Family** without their love, encouragement and support, I would not have been able to achieve my goals. Thank you from the bottom of my heart. **I love you all.**

I would like to thank my **all Friends** for their help and offer me advice.

Thanks for **everything** that helped me get to this day.

Thanks for Faculty of Engineering.

Thanks **Egypt.**

Table of Contents

ACKNOWLEDGMENTS	I
TABLE OF CONTENTS	III
LIST OF TABLES	V
LIST OF FIGURES	VI
LIST OF ACRONYMS	VIII
ABSTRACT	IX
CHAPTER 1: INTRODUCTION	1
1.1. PROBLEM STATEMENT	1
1.2. MOTIVATION	2
1.3. CONTRIBUTIONS	3
1.4. ORGANIZATION OF THE THESIS	3
CHAPTER 2 : BACKGROUND	4
2.1. OVERVIEW OF DATA MINING	4
2.1.1. Data Mining Techniques	5
2.1.1.1. Predictive Mode	
2.1.1.2. Descriptive ModeL	
2.2. ASSOCIATION RULE MINING	
2.2.1. Association Rule Mining Algorithms	
2.2.1.1. Apriori Algorithm	
2.3. BIG DATA	
2.3.1.1 Spark	
CHAPTER 3 : RELATED WORK	
3.1. INTRODUCTION	
3.2. APRIORI ALGORITHM BASED ON MAPREDUCE	
3.2.1. The proposed Algorithm by Zeng L et.al.[15]	21

3.2.2. The proposed Algorithm by Brijendra Singh et.al.[16]	22
3.2.2.1 BloomFilter	22
3.3. APRIORI ALGORITHM BASED ON SPARK	24
3.3.1. The proposed Algorithm by Hongjian Qiu et.al.[18]	24
3.3.2. The proposed Algorithm by Rathee et. al.[19]	
3.3.3. The proposed Algorithm by Krishan Kumar Sethi et. al.[20]	
3.4. SUMMARY	
CHAPTER 4: THE PROPOSED APPROACH	32
4.1. APRIORI PERFORMANCE WITH INCREASING TRANSACTIONS	32
4.2. OVERVIEW OF CUCKOO FILTER STRUCTURE	34
4.2.1. The Cuckoo Filter Algorithms	35
4.2.1.1. Insert Operation	36
4.2.1.2. Loockup Operation	38
4.2.1.3. Delete Operation	38
4.2.2. False Positive Rate Of a cuckoo Filter	39
4.2.2.1. Optimal bucket size(b)	39
4.2.3. Space Complixty	42
4.2.4. Time Complixty	42
4.3. THE ASCF ALGORITHM	43
CHAPTER 5 : EXPERIMENTAL RESULTS	58
5.1. EXPERIAMENT SETUP	58
5.1.1. Datasets	58
5.2. PERFORMANCE ANALYSIS	59
5.3. DISCUSSION	
CHAPTER 6 : CONCLUSION AND FUTURE WORK	67
6.1. CONCLUSION	67
6.2. FUTURE WORK	68
REFERENCES	69
APPENDIX A: THE ASCF CODE	72

List of Tables

Table I.1: List of AcronymsVIII
Table 2.1: Calculated Support8
Table 2.2: Calculated Confidence8
Table 2.3: Example of Horizontal Layout Database9
Table 2.4: Example of Vertical Layout Database10
Table 2.5: Transaciton Dataset D
Table 2.6: 1-Frequent Itemsets
Table 4.1: Description of Diabetic Patient Datasets
Table 4.2: The Time Complexity of Cuckoo Filter
Table 5.1: Dataset and Their Attributes59
Table 5.2: Execution Times per Iteration on T10I4D100K Dataset for ASCF, HFIM and
YAFIM61
Table 5.3: Execution Times per Iteration on Chess Dataset for ASCF, HFIM and YAFIM
61
Table 5.4: Execution Times per Iteration on Retail Dataset for ASCF, HFIM and YAFIM
62
Table 5.5: Total Execution Times for ASCF, HFIM and YAFIM62
Table 5.6: The Number of Transactions, Frequent Items, Frequent Itemsets for
T10I4D100K Dataset64
Table 5.7: The Number of Transactions, Frequent Items, Frequent Itemsets for Retail
Dataset64
Table 5.8: The Number of Transactions, Frequent Items, Frequent Itemsets for Chess
Dataset

List of Figures

Figure 2.1: Data Mining Technique	5
Figure 2.2: Apriori Algorithm	
Figure 2.3: Transactional Dataset	
Figure 2.4: Generating 1-Frequent Itemsets	12
Figure 2.5: Generating 2- Frequent itemsets	
Figure 2.6: Generating 3- Frequent itemsets	
Figure 2.7: Construct FP-Tree	
Figure 2.8: Example of Transformation	
Figure 2.9: Example of Action	18
Figure 2.10: Transformed RDD are Computed Iazily	19
Figure 2.11: Spark RDD Lineage Graph	
Figure 2.12: Working Model of Spark Framework	20
Figure 3.1: Bloom Filter	23
Figure 3.2: YAFIM Lineage Graph for the RDDs in Phase one	24
Figure 3.3: YAFIM Lineage Graph for the RDDs in Phase two	25
Figure 3.4: R-Apriori Lineage Graph for the RDDs in Phase one	26
Figure 3.5: R-Apriori Lineage Graph for the RDDs in Phase two	27
Figure 3.6: R-Apriori Lineage Graph for the RDDs in Phase three	28
Figure 3.7: HFIM Lineage Graph for the RDDs in Phase one	29
Figure 3.8: HFIM Lineage Graph for the RDDs in Phase two	30
Figure 4.1: Execution Time for Hospital Data	
Figure 4.2: Execution Time for Diabetic Patient Data	34
Figure 4.3: Insert the Items {a, d, e} into Cuckoo Filter	35
Figure 4.4: Algorithm 1 Insert (x) into Cuckoo Filter	37
Figure 4.5: Algorithm 2 Lookup(x) from Cuckoo Filter	38
Figure 4.6: Algorithm 3 Delete(x) from Cuckoo Filter	39
Figure 4.7: Load Factor Achieved by Using a vary f-bit fingerprint with Different S	Size
of Buckets	40
Figure 4.8: Amortized Bits per Item Measured False Positive Rate with Different	
Bucket Size	41
Figure 4.9: ASCF Lineage Graph of RDDs in Phase one	44

Figure 4.10: ASCF Algorithm Phase one	45
Figure 4.11: The Flow Chart of ASCF in Phase two	
Figure 4.12: ASCF Algorithm Phase two	49
Figure 4.13: ASCF Lineage Graph of RDDs in Phase two	51
Figure 4.14: Generating 1-Frequent Itemsets	52
Figure 4.15: Generate (Combination of 2 items, 1) from Partition 1	53
Figure 4.16: Generate (Combination of 2 items, 1) from Partition 2	53
Figure 4.17: Generate (Combination of 2 items, 1) from Partition 2	54
Figure 4.18: Generate 2-Frequent Itemsets	54
Figure 4.19: Modify Cuckoo Filter	55
Figure 4.20: Generate (Combination of 3 items, 1) from Partition 1	56
Figure 4.21: Generate (Combination of 3 items, 1) from Partition 2	56
Figure 4.22: Generate (Combination of 3 items, 1) from Partition 3	57
Figure 4.23: Generate 3-Frequent Itemsets	57
Figure 5.1: Execution Times for ASCF, HFIM and YAFIM on Three Datasets	60

List of Acronyms

Acronym	Definition
ARM	Association Rule Mining
ASCF	Apriori Algorithm on Spark based on Cuckoo Filter structure
HDFS	Hadoop Distributed File System
HFIM	a Spark-based Hybrid Frequent Itemset Mining algorithm for big data processing
KDD	Knowledge Discovery in Databases
LHS	Left Hand Side
R-Apriori	An Efficient Apriori based Algorithm on Spark
RDD	Resilient Distributed Dataset
RHS	Right Hand Side
YAFIM	(Yet Another Frequent Itemset Mining) A Parallel Frequent Itemset Mining Algorithm with Spark

Abstract

Data mining is the process that is used for extracting interesting patterns from large amount of data using a variety of techniques. One of the techniques that help to discover important relations between variables in large dataset is Association rule mining (ARM). It is used to identify strong rules discovered in databases. To build the association rules for all itemsets this requires large memory and processing resources. So, only frequent itemsets are considered to reduce the number of itemsets. There are many frequent itemset mining algorithms; the most popular algorithm is the Apriori algorithm. The Apriori algorithm is iterative and works sequentially, its normal execution is on a single machine. Nowadays, there is a great explosion of data; the rapid growth of data necessitates running the data intensive algorithms in parallel distributed environment to achieve convenient performance. The major drawbacks of the Apriori algorithm concerning its computational complexity make the algorithm inefficient to use when the data size is getting larger. In this thesis, Apriori Algorithm on Spark based on Cuckoo Filter structure (ASCF) is introduced. The ASCF algorithm solves the inherent drawbacks in the original Apriori algorithm. It succeeds in removing the candidate generation step from Apriori algorithm to reduce computational complexity and avoid costly comparisons, and uses cuckoo filter to further enhance the performance. The proposed algorithm is implemented on spark inmemory processing distributed environment to reduce processing time. It offers great improvement in performance over other previous approaches of Apriori algorithm implementation based on spark. The ASCF on a cluster of 4 nodes achieves a time of only 5.8% of the competing approach on the Retail dataset with minimum support of 0.75%, 25.6% on Chess dataset with minimum support of 85% and 37.3% on T10I4D100K with minimum support of 0.25%.

Chapter 1: Introduction

We live in the data age, this data is being generated by everything around us at all times: from social sites, sensors, search engine, medical reports, etc. This huge amount of data is known as 'Big Data'. It is as a collection of large, diverse, complex dataset that comes at high speed. There is an urgent need to assist humans in extracting useful information (knowledge) from this data. It requires new architecture, techniques, algorithms, and analytic methods to manage it and extract value and hidden knowledge from it. Big Data mining is the process of discovering patterns from large datasets that have those mentioned properties.

1.1. Problem Statement

The problem can be formally stated as:

Let D be a set of transactions in large transactional data where each transaction T is a set of items such that $T \subseteq I$, where $I = \{i1, i2, \dots, im\}$ is a set of items. Each transaction has a unique identifier TID. Let X and Y be two different sets of items, where $X \subseteq I$, $Y \subseteq I$ are subset of I and $X \cap Y = \emptyset$. If $X \subseteq T$ then the transaction T is said to contain X.

Each itemset appears a number of times in D so there are a percentage of transactions that contain all items in X, which is called support(X). The itemsets X and Y appear together in D if there are S percentage transactions that contain X \cup Y, where S is the $support(X \cup Y)$, we say there is association between itemsets X and Y.

The problem is discovering important association rules between variables (items) in large transactional dataset. These rules are called strong rules if they have support and confidence greater than or equal predefined minimum support and confidence thresholds.

Mining association rule task consist of two steps:

- Finding all frequent itemsets which have support greater than or equal a predefined minimum support
- Generating confident rules from the frequent itemsets discovered.

The performance of an algorithm for mining association rules is determined by the first step, while the second step is straightforward. Most of the algorithms are inefficient and take a lot of time when working with large datasets because they are designed to work with a single machine in an iterative manner with no natural way for distribution. So, new technologies are required to store and process large data sets in a distributed computing environment like Hadoop Mapreduce and Spark.

1.2. Motivation

The Apriori algorithm is iterative and works sequentially. Its normal execution is on a single machine whose speed leads to severe performance degradation when working with big data so multiple machines and a parallel algorithm are needed where issues like data replication and synchronization must be addressed. The major drawbacks of the Apriori algorithm concerning its computational complexity make the algorithm inefficient to use when the data size get larger. It generates candidate sets having all possible pairs from frequent itemsets that are generated from previous iteration and compares each pair with every transaction record in every iteration.

In this thesis, we apply Apriori Algorithm on Spark based on Cuckoo Filter structure (ASCF) to solve these problems. We propose a parallel Apriori implementation on Spark distributed environment to reduce processing time, and use cuckoo filter structure to further enhance the performance. We focus on the step that is the most time consuming in the whole algorithm, which is generating candidate sets in the second phase and solve this problem by eliminating it to reduce computational complexity and avoid costly comparisons.