

## **Comparison of three association rule mining algorithms using CloudSim**

A Dissertation Proposal Submitted By

**Mandeep Kaur**

**(11301803)**

To

**Department of Computer Science & Engineering**

In partial fulfillment of the Requirement for the Award of the Degree of

**Master of Technology in Computer Science & Engineering**

Under the guidance of

**Mr. Gaurav Kumar Tak**

**(May, 2015)**

School of: Computer Science and Engineering

**DISSERTATION TOPIC APPROVAL PERFORMA**

Name of the Student: Mandeep Kaur

Registration No: 11301803

Batch: 2013-2015

Roll No. B54.

Session: 2014-15

Parent Section: K2305

**Details of Supervisor:**

Designation: Assistant Professor

Name: Gaurav Kumar Tak

Qualification: M.Tech

U.ID 15746

Research Experience: ---

**SPECIALIZATION AREA: Data Mining (pick from list of provided specialization areas by DAA)**

**PROPOSED TOPICS**

1. Associate Rule Mining with FP-Growth in Cloud SIM.
2. Apriori Algorithm.
3. Data Mining techniques on swarm intelligence algorithm

*M Gaurav Kumar*  
*15746*  
Signature of  
Supervisor

**PAC Remarks:**

*Topic one is approved. Paper publication is also completed*

**APPROVAL OF PAC CHAIRPERSON:**

Signature: *(Signature)*

Date: *26/9/14*

\*Supervisor should finally encircle one topic out of three proposed topics and put up for approval before Project Approval Committee (PAC)

\*Original copy of this format after PAC approval will be retained by the student and must be attached in the Project/Dissertation final report.

\*One copy to be submitted to Supervisor.

## ABSTRACT

Association rule mining is one of the methods of data mining techniques which is used to find out the interesting, familiar correlations and find out association between set of items in the transactional database. These algorithms which are discussed in this paper namely Apriori algorithm, FP growth algorithm, Improved Apriori algorithm compared on the basis of different type of dataset, support count and some other factors. These three algorithms are enhanced using CloudSim examples which leads to a collaborated approach of cloud and mining. Both made collectively a cloud mining approach in which these three algorithms are implemented. Cloud Computing and Data Mining both are the growing approaches in the present scenario of information technology. Data mining method is used for finding knowledge for the rare data in addition to Cloud computing give us secure as well as stretchy communications which tells the whole thing as a service. Through data mining in addition to cloud computing integration provides quickness and rapid admittance in the direction of technology. The consequence of such mixing is supposed to be in the direction of powerful and capacitive stage to motivate the talented contract by means of the growing manufacture of data, in order to create the circumstances in favor of the well-organized mining of big quantity of data consisting of variety of data warehouses by means of the aspire of creating supportive knowledge or manufacturing of fresh information. The idea to make big quantity of data passes out the abilities to find and use precious knowledge of data. Data mining have been a winning instrument to examine data on or after dissimilar angles and receiving helpful knowledge from data. It will be able to well assist in finding standards, categorization of information, labeling of information as well as in the direction to discover patterns similar of the dataset.

Cloud communications are able to be used efficiently for comprehensive environments in addition to challenging operations through data so that the process of data mining can be carried out in distinctive manner.

## **CERTIFICATE**

This is to certify that **Mandeep Kaur** has completed M.TECH dissertation title “**Comparison of three association rule mining algorithms using CloudSim**” under my guidance and Supervision. To the best of my knowledge, the present work is his original investigation and study. No part of dissertation has ever been submitted for the any degree and diploma. This dissertation is fit for the submission and partial fulfilment of the condition and for the award of M.TECH computer science and engineering.

Signature of the advisor

Name- Gaurav Kumar Tak

UID- 15746

## **ACKNOWLEDGEMENT**

I would like to place on record my deep sense of thankfulness to Mr. Gaurav Kumar Tak of Dept. of computer science and engineering, Lovely Professional University Jalandhar, India for his substantial supervision, assistance and valuable proposals. I express my sincere thankfulness to Asstt. Prof. Dalwinder Singh, Dept. of Computer science and Engineering, LPU, Jalandhar, India, for his inspiring supervision, nonstop reinforcement and administration throughout the development of existent work. For any errors or inadequacies that may remain in this work, of course, the responsibility is entirely my own.

MANDEEP KAUR

## **DECLARATION**

I hereby state that the dissertation proposal entitled, “comparison of three association rule mining using CloudSim” , Mandeep Kaur submitted for the M.Tech Degree is completely my unique effort and all thoughts plus references have been suitably recognized. It does not hold any effort for the reward of any other grade or certificate.

**Date:** 05-05-2015

**Investigator**

**Reg. No.:** 11301803

## TABLE OF CONTENTS

Chapter number	Page number
1.Introduction 1	
1.1 Data Mining .....	1
1.2 Data Mining Techniques.....	2
1.2.1 Classification.....	3
1.2.2 Clustering .....	3
1.3 Association Rule Mining .....	3
1.4 Apriori Algorithm .....	5
1.5 FP-Growth Algorithm.....	7
1.6 Cloud computing.....	7
1.6.1 On-Demand Self-Service .....	8
1.6.2 Broad Network Access.....	8
1.6.3 Resource Pooling.....	9
1.6.4 Rapid Elasticity .....	10
1.6.5 Metered Service.....	10
1.7 Data Mining with Cloud Computing .....	10
2. Review of literature .....	13
3. Scope of the study.....	23
4. Objectives of study .....	24
5. Research methodology.....	26
5.1 Flowchart of work.....	26
5.2 Tool Used.....	27
5.3 Apriori algorithm .....	28
5.4 FP growth algorithm .....	29
5.5 Improved Apriori algorithm.....	30
6. Results and discussions .....	34
6.1 Experimental work.....	34
6.2 Data analysis and interpretation.....	38
7. Conclusion and future scope.....	42
8. References.....	43
9. Appendix.....	47

## TABLE OF FIGURES

Figure 1.1 Data mining as KDD.....	2
Figure 1.2 Data mining techniques.....	2
Figure 1.3 Cloud computing characteristics.....	7
Figure 1.4 Cloud services.....	11
Figure 5.1 Flowchart of work.....	26
Figure 5.2 FP tree.....	30
Figure 6.1 Register user page.....	34
Figure 6.2 login user page.....	35
Figure 6.3 Upload and fetch dataset.....	36
Figure 6.4 Selection of user for data mining.....	37
Figure 6.5 Results.....	38
Figure 6.6 Analysis of time and memory between vertical Apriori, horizontal Apriori and FP growth algorithm.....	39
Figure 6.7 Results.....	40
Figure 6.8 Analysis of time and memory between vertical Apriori, horizontal Apriori and FP growth algorithm.....	41



## LIST OF TABLES

Table 5.1: Generation of candidate tables from database .....	29
Table 5.2: Database for FP growth example .....	30
Table 5.3: Database of horizontal format .....	32
Table 5.4: Conversion from horizontal to vertical database .....	32



With the advancement in Information Technology, Now a days the size of the databases created by the organizations has been widely increased likewise due to the present of the advancement in the data capturing technologies. It is possible to store large amount of organizational data and scans it. This kind of data includes transaction records in supermarkets, banks, stock markets, and telephone companies. Data however is not the same as information, it should be analyzed and extracted before it is useful .It is a challenge to handle such explosive growth of amount of data, and to find techniques to find out useful knowledge out of such a massive history of transactional data. Data mining has been emerged as a new research area to meet this challenge. The knowledge that is needed through the help of data mining can be applied into various applications like business, retail and market, engineering design and scientific area.

#### 1.1 Data Mining

Data mining can also be known as knowledge discovery in databases as for the detail way, data mining has been recognized as very important research area to efficiently extract useful information from large databases. Knowledge Discovery in Data mining has become one of the most popular and important research areas in the database community. In recent years, data mining has been used in every kinds of areas of science and engineering, also in genetics, bioinformatics, medicine engineering. Also people from business area find more and more applications for data mining for the research, most applications are found in retail, insurance, security, telecommunication. If we want to go in detail view of the data mining then we are always advised to refer the knowledge discovery in broader way. This is most exciting data mining technique in now a days. For referring the research more popular and useful we select the data mining.

The KDD process includes an iterative sequence method:

- **Selection:** In the KDD process main concern is firstly to select the data needed for data mining from the different and heterogeneous data sources.
- **Preprocessing:** includes missing data or faulty data. To find these kind of noise there are many different activities involved. Data which contains error may be

corrected or removed. In the preprocessing of data we remove the noise or outliers, collecting important information to model for noise.

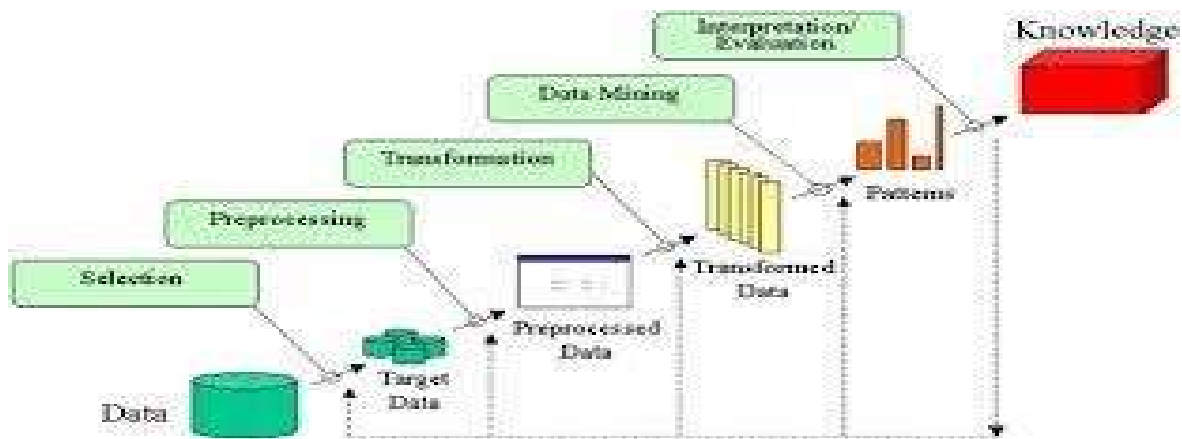


Figure 1.1: Data Mining as KDD

- **Transformation:** is used to translate the data into the common format for preprocess of data. Some of data may be converted into more usable format. Data reduction method may be used to decrease the number of imaginable data values being considered.
- **Data Mining:** is used for making the useful results from the data. Using data mining we can find the hidden patterns inside the huge amount of data.
- **Interpretation/Evaluation:** is used to presented the data mining outcomes to the users which are significant because of the usefulness of the result is depend on it. There we can use numerous visualization plus GUI approaches in this phase.

## 1.2 Data Mining Techniques

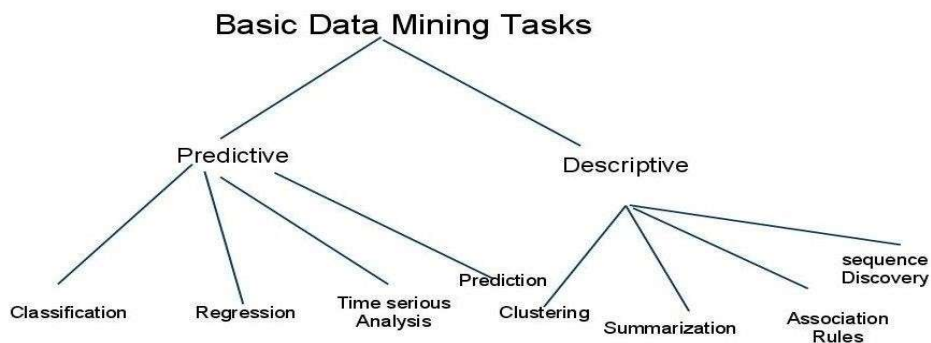


Figure 1.2: Data Mining techniques

### **1.2.1 Classification**

Classification is one of the common data mining techniques. In which a set of pre-classified examples is used to create a model which can classify the population of records. In this we discover the predictive learning function which is used to classify a data item into one of several predefined classes. Credit risk applications are well suited for this type of analysis. Data classification process includes the learning and in learning/training data both are analyzed by classification algorithm.

### **1.2.2 Clustering**

Clustering is also an important and well known data mining technique. Clustering is defined as the process in which we form the clusters which means identification of similar classes of objects. It is a common descriptive task in which we form the finite sets according to their common identity. If we use clustering techniques then we can discover the overall distribution pattern.

### **1.3 Association Rule Mining**

Association rules are one of the most common and important techniques of data mining that is used to find out hidden and interesting patterns or associations among the data items stored in the database scan. The invention of interesting association relations among huge amount of records can help in decision making process, such as catalog design, marketing etc.

The mining technique of association rule consists of two sub steps (1) judgment all common item sets that occur more than a minimum support threshold (2) generate association rules by these common item sets. The first sub step shows a significant role in association rules mining plus has a time complexity and is more costly. Association rule mining is mainly used in market basket analysis and retail data analysis. In market basket analysis we refer to the different buying habits of customers and study them to find associations among items those are purchased by other customers. Items which are frequently purchased together by customers are the most common frequent sets. The methods for discovering association rules from the data that dedicated on discovering connections between items telling some characteristic of a retail environment, frequently purchasing behavior of customers for determining items that customers buy together all the time. All rules of this type throw light on particular local pattern. The association rules can be straightforwardly understood.

All association rule is in the form of “ $A \Rightarrow B$ ” always where A and B are disjoint and unique item sets ,

$$\text{i.e., } A \cap B = \emptyset$$

The percentage of an association rule can be calculated in terms of its confidence and support. Association rule mining is used to discover out association rules that fulfill the predefined minimum support threshold plus confidence from a record. The support of an item in the database is the percentage of transactions in which that item occurs frequently. The confidence is used to find out the percentage of the rule and the ratio of the number of transactions that contain A or B to the number of transactions that contain A.

$$\text{Confidence } (A \Rightarrow B) = P(B | A) = \frac{\text{support-count}(A \cup B)}{\text{support-count}(A)}$$

The rules which are finding by the association rule mining that satisfy both the user specified minimum support threshold plus confidence can be said to be the Strong Association rules. The determination is to examine for the relationships among items of database. That's why, association rules can help decision creators to find out the likely items that are likely to be bought by customers composed always. Association rule mining is an significant research area of data mining; its task is to discover all subsets of items which repeatedly occur, and the connection between them. Association rule mining has two major steps: the establishment of common itemsets and the establishment of rules. Association rule mining is defined as to check out rules of association which fulfill the predefined least support after a transaction database. Both support and confidence are most important terminologies used in finding the association rules. When we are going with association rule problem it is sub rotten into two sub problems. First sub problem is to discover out those itemsets who exceed the occurrence of predefined support in the database: known as frequent itemsets. The second sub problem is to produce association rules from those common itemsets under the confidence given. Main motive of association rule is to find out correlation among the different transactions , it make easier for taking decisions and to use the process efficiently. There are so many algorithms which are used to mine frequent itemsets. Normally, all association rule mining algorithm contains the subsequent steps:

- i. The group of candidate J-itemsets is produced by 1-extensions of the large (J -1) itemsets produced in the preceding repetition.
- ii. Support for the candidate J-itemsets are produced by a pass directly above the transaction database.
- iii. Itemsets which are not obeying the minimum amount of support are not necessary and the residual itemsets are named oversized J-itemset.

In elaborated way an association rule is an manifestation like  $(P \Rightarrow Q)$  where (P) and (Q) are set of items. The definition of these kind of rules is reasonably remarkable. Given a transaction database (X) where each transaction ( $T' \in X$ ) is set of items.  $(P \Rightarrow Q)$  signifies that whether that transaction ( $T'$ ) contains (P) than ( $T'$ ) probably contains (Q) also. The probability can be explained as the proportion of transactions covering (Q) in contrast To (P) with concern to global number of transaction holding (P). That is, confidence can be also known as the conditional probability  $P(Q)$ . The awareness of mining association rules invents from the investigation of market basket data in which rules like “ A somebody who purchase ( $Y_1$ ) with ( $Y_2$ ) will also purchase item (Y) with possibility c% “ are originate. Their straight applicability to commercial problem composed with their inborn recognize aptitude for non-data mining methodologies.

There are large amount of algorithms which have main point to mine common happening itemsets. Certain of them are very fine identified which run a entire new period in data mining field. They come up with mining common itemsets and association rules likely. The algorithms deviate essentially in what way the candidate sets are produced plus in what way the supports for the candidate sets are scheduled.

Various algorithms like partition algorithm, Apriori algorithm, dynamic item set counting algorithm, FP tree growth algorithm have been developed to find out the frequent item set from the transaction database.

#### **1.4 Apriori Algorithm**

The Apriori algorithm is the most widely used association rule mining algorithm. In this firstly generate frequent itemsets and generate association rules. The basic idea behind this is scanning the database again and again till the end we did not find the interesting

patterns or rules. This principle applied on the database many times. Apriori algorithm uses the tree like structure to count candidate item sets effectively. Apriori algorithm follows the breadth first search. In this we generate the (i) candidate itemset from (i-1) itemsets. Candidate itemsets having (i) items can be formed using joining frequent itemsets having

( i-1) items, and by removing all other subsets which are not frequent. Mainly the algorithm constitute of the two main steps i.e itemsets generation and association rule mining. Again frequent itemsets generation is of two-step process:

Candidate itemsets generation in which all possible combination of items those are frequent itemsets.

Frequent itemsets generation which support for all candidate item sets are generated and those itemsets which have greater minimum support than the user specified are referred as the frequent itemsets. Apriori algorithm scan the database over the multiple times till the end find out the frequent itemset at last. Apriori algorithm has been extremely commonly used for mining of common item sets and to find out associations. The main dissimilarity in Apriori is the fewer candidate itemsets it produces for analysis in every database go by. The exploration for association guidelines is showed by 2 key points: support along with that confidence. Apriori algorithm results constantly return an rules of association if its support plus confidence morals are higher than user specified values. The outcome is organized by confidence. If numerous rules have the related confidence then they all are supported by support and are the best rules. Thus Apriori favoritism extra best confident rules and describes these rules as extra remarkable .Apriori algorithm also uses a penetrating method called breadth first search method. It calculates the support of itemsets. It has a candidate generation fragment which makes use of the downward closure property of the support count. There are basic two chief steps of the Apriori algorithm are namely the join and prune steps.

(a) The join step is basically used to construct fresh candidate sets. A candidate itemset is a itemset which can be either frequent or infrequent with regard to the support threshold given. Highest level candidate itemsets ( $C_i$ ) are produced by joining last level occurring itemsets are ( $L_{i-1}$ ) with new one.



(b) The prune step used in checking out and to decrease out candidate item-sets whose last levels are not common. This is mainly founded on the anti-monotonic property as a outcome of which each subset of a common item set is also common. Thus a candidate item set which is collection of one or more uncommon item sets of a last level is cleaned(pruned) from the method of common itemset and association rule mining. Apriori Algorithm.

### **1.5 FP-Growth Algorithm**

FP-Growth algorithm uses important feature like section development technique which is used to evade the huge amount of scans on the transactional database. This algorithm involves only binary scans. FP-Growth algorithm's initial step is to calculate a list of common items which are organized by frequency in downward instruction (F-List) through its initial transactional database scan. In instant scan, the innovative transaction database is compacted or turned into a highly condensed FP-tree. Afterward the manufacture of FP-tree from that, FP-Growth mechanism in a divide plus conquer mode and the FP growth algorithm accomplishes mining on FP-tree.

Apriori algorithm always needs  $(i+1)$  scans, where  $(i)$  is the length of the longest outline, we can use FP growth technique to condense the amount of scans of the entire transactional database D to find the frequent itemsets using only two scans of database.

### **1.6 Cloud computing**

Cloud computing can be defined on the basis of their essential characteristics like as below figure showing:

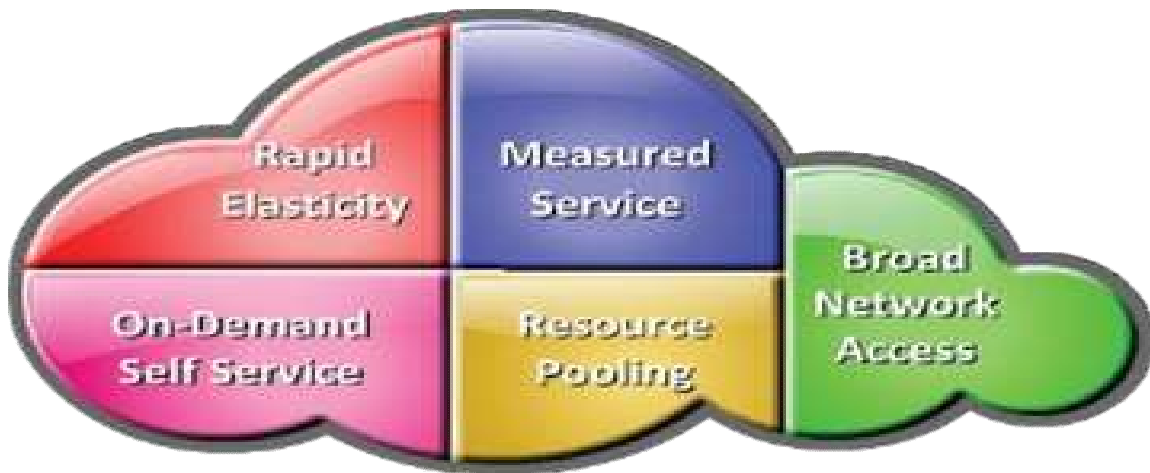


Figure 1.3: cloud computing characteristics

### 1.6.1 On-Demand Self-Service

- This service ensures customers to get the resources required when they need, without any involvement of human resources.
- It provide customers the facility named always ready to use service, which provided facility of any time use from the catalog.
- It does not restrict the users to use the self-service interface
  - ▶ Self-service interface provided always by user-friendly environment.
- The on-demand self-service features of Cloud Computing means that a customer can customize the Cloud services as mandatory, without any human resources with the Cloud service provider. With usage of the self-service interface, customers can accept Cloud services by demanding for the essential IT capitals from the service catalog list.

### 1.6.2 Broad Network Access

- Cloud services of broad network access are retrieved via the system, commonly from the Internet and from a wide-range of client stages such as:
  - ▶ Ordinary PC's
  - ▶ New PC's
  - ▶ Cellphones

► Thin Client access

- Removes the requirement for retrieving a specific client platform to use the services.
- Allows retrieving the services from everywhere across the globe.
- Users have to connect the software on their pc's in demand to practice this software request. It is not desirable to admit this software if the operator is absent from the PC's where the software is connected. Nowadays, abundant of the recycled softwares can be retrieved above the Internet. For instance, Google Docs, a Web-based manuscript maker and copyreader permits operators to admission

and edit papers from any device with an Internet connection, removing the essential to have admission to a specific client platform to control papers.

### **1.6.3 Resource Pooling**

- IT capitals (calculate, storing, system) are joint to attend manifold customers, Centered on multi-tenant model.
- Customers have no information of the particular site of the capitals delivered to them.
- Capitals are animatedly allocated and reallocated founded on the customer request.
- A Cloud necessity has a big and stretchy supply pool to encounter the customer's essentials, to deliver the markets of gauge, and to encounter service-level necessities. The capitals from the pool are animatedly allocated to manifold customers founded on a multi-tenant model. Multitenancy mentions to an design and scheme by which manifold self-governing tenants are repaired by a solitary set of capitals. In a Cloud tenant could be a operator, a user collection, or an group/corporation. Multitenancy allows calculate, storing, and system resources to be common among manifold clients. Virtualization delivers habits for allowing multitenancy in Cloud. For instance, manifold VMs after dissimilar customers can run concurrently on the identical server with hypervisor provision.

- There is a wisdom of location individuality, in that the customer usually has no knowledge of nearby precise location of the capitals delivered .

#### **1.6.4 Rapid Elasticity**

- Capability is to measure IT capitals quickly, as compulsory, to achieve the fluctuating requirements minus disturbance of service. Capitals can be together climbed up and climbed down animatedly.
- To the customer, the Cloud seems to be immeasurable. Customers can jump with negligible calculating control and can enlarge their surroundings to any extent.
- Rapid elasticity mentions to the capability of the Cloud to enlarge or decrease assigned IT resources rapidly and resourcefully. This distribution might be done mechanically to any service disruption. Customers will take benefit of the Cloud when they have huge variation in their IT supply usage. For instance, the association may be compulsory to dual the amount of Web and application attendants for the whole period of exact job. They would not need to wage the capital expenditure of having idle servers on the bottom furthestmost of the period and that too would need to discharge these servers capitals after the job is accomplished. The Cloud allows to produce and analyze these capitals animatedly and lets the establishments to wage these on a practice basis.

#### **1.6.5 Metered Service**

- Customers are payable on the metered procedure of Cloud capitals.
  - ▶ Price experienced on a pay-per-use source.
  - ▶ Rating/publicizing prototypical is knotted up with the obligatory facility levels.
  - ▶ Source usage is watched and described, which delivers transparency for chargeback to together Cloud service supplier and customer nearby the utilized service.

#### **1.7 Data Mining with Cloud Computing**

Cloud Computing is a basic business model. Which divides the computing responsibilities to the supply pool which consist of a huge number of workstations, by

which various application methods can achieve storage space plus software facilities on request. Cloud computing is basically used to deliver the computing plus storage volume as a facility to addressees. The main importance lies when we get to know that if we use the cloud then we have not to worry to run the software from our computer because it is stored on the server. The more important feature is that if a computer crashes, the other users can still use the software.

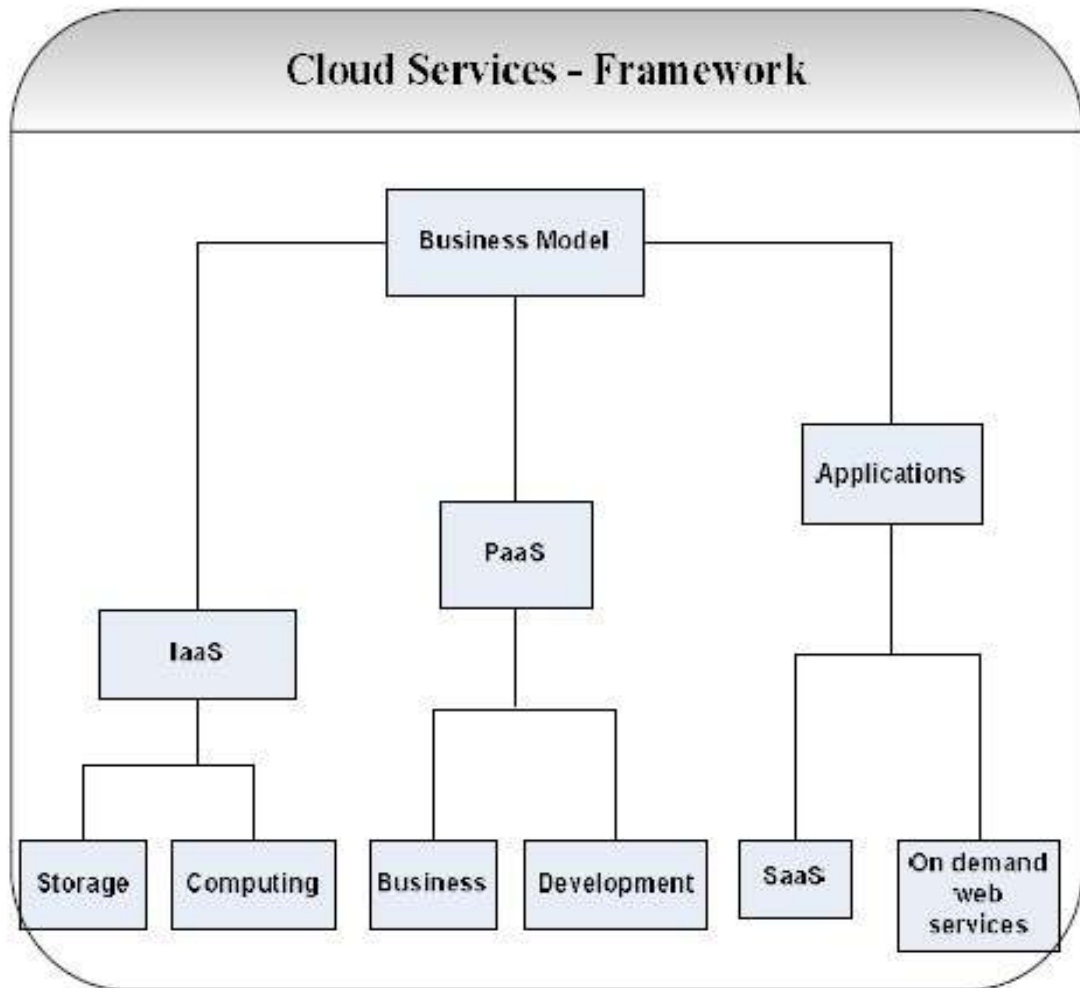


Figure 1.4: Cloud Services

Now a days data mining and cloud computing both permits business to consolidate the supervision of software with the guarantee of well-organized, consistent services for their customers. This collaborated approach provides tools that can hold huge quantity of data and provide cheap and efficient solution. Cloud mining is the revolutionary term now days for data mining. There is huge amount of data which are very difficult to mine and find out meaningful knowledge the in terms of computational resources. Using cloud computing in the data mining we can access the data more easily and efficiently as well

as the cost effective computational resources. Data mining in the cloud computing permits many administrations to centerline the supervision of software and data with the guarantee of consistent and confident facilities for their customers. Cloud computing is refereed to hardware and software supplied as facilities above the Internet, and in cloud mining softwares can deliver facilities in this way.

The most important properties of data mining tools provided through the cloud are:

- The customers simply give money for the data mining apparatuses that he or she needs.
- The customers does not retain a hardware organization.

This methodology also decreases the issues that preserve some minor corporations from promoting of the data mining tools. The collaboration of cloud computing conveys new innovative thoughts for data mining and it upsurges the measure of handling data at huge level.

**Qihua Lan et al** shown in their paper an innovative algorithm approach entitled APFT , which syndicates the Apriori algorithm with FP-tree structure and they both collaborately formed FP-growth algorithm. The benefit of Apriori with FP tree is in this algorithm its not essential to produce unconfirmed pattern bases plus sub- unconfirmed pattern tree. At last the outcomes of the experimentations display that it mechanisms almost as quick as FP-growth also quicker than Apriori (Qihua Lan et al, 2014).

**Lingjuan Li, Min Zhang** focuses on the approach of association rules mining in cloud computing surroundings and introduces Hadoop, Apriori algorithm plus parallel rule mining association algorithm and MapReduce programming model. In this paper a comparable rule mining association approach designing which adjusted to the cloud computing surroundings. The data set partition technique is planned in this paper are also enriched with Apriori algorithm; which is executed the enhanced algorithm established on MapReduce programming model on the Hadoop platform. At the end outcomes demonstration that the approach intended in this paper can achieve advanced efficiency when we are doing common itemset mining in cloud computing surroundings(Lingjuan Li, Min Zhang, 2014).

**Shruti Mishra et al** have fuzzified the association numerous common pattern mining methods and fuzzified original dataset which have been applied to invent meaningful frequent patterns. The paper has put a clear judgement of the common pattern mining methods in the fuzzified and original data in terms of attributes like runtime of the algorithm with the number of common patterns produced. At the end the resultset showed that the fuzzified set is capable of inventing a huge number of common patterns and have an enhanced run time duration(Shruti Mishra et al, 2013).

**T.V.Mahendra et al** discussed of an algorithm to mine the data by the cloud structure consuming sector/sphere structure through association rules. Well known Data mining is defined as the process of processing data from dissimilar aspects and longwinded it into meaningful material. Association rule mining is one of the supreme essential

prospect in data mining. Rules of Association are dependent on each other which formed happening of an item centered on happenings of further items. Apriori is referred to familiar procedure to mine rules of association. Cloud can be intended as an organization which delivers possessions of facility above the internet. A Cloud service can be a compartment cloud that delivers block or file centered storage facility or it can be a figured out cloud that provides computational services. In addition this paper has focused on the implementation and design of sector and sphere storage cloud. Sector is referred to the scattered file scheme, however sphere is comparable in-storage data dispensation structure that can be recycled to development data stored in sector(T.V.Mahendra et al, 2013).

**Tipawan Silwattananusarn et al** discover the uses of data mining methods which have been established to grasp up information organization procedure. The periodical articles indexed in Science Direct Database from 2007 to 2012 are categorized. The argument on the discoveries is allocated into 4 areas: (a) information supply; (b) information categories and/or information datasets; (c) information mining responsibilities (d) information mining methods that uses with information association. The object in transitory designates the explanation of data mining along with data mining workings. Afterward the information management basis and chief information management apparatuses comprised in information management phase are defined. To complete, the requests of data mining methods in the procedure of information management are discussed(Tipawan Silwattananusarn et al, 2012).

**K.Ganeshkumar et al** presents a set of encryption approaches for transactional databases which are applicable for contract out rules of association. Starting as of an easy one-to-one replacement cipher, the paper developed a complicated conversion algorithm that inserts non-deterministic info to the deterministic outcomes of a one-to-n item representing arrangement .The outcomes display that encryption method is very strong to occurrences as conflicting to modest one to one ciphers, which can be simply cracked with the support of upbringing information. The experimentations presented that the encryption price is reasonable and much inferior than the mining price(K.Ganeshkumar et al, 2012).

**Ankit Bhardwaj1 et al** focuses on the data mining and the present trends connected with it. It presents an review of data mining system and clarifies how knowledge discovery



and data mining in databases are linked mutually to all other and to correlated fields. The numerous data mining techniques are introduced by the unusual researchers. These techniques are used to perform classification, to accomplish clustering, to discover attractive patterns. In their potential work, the data mining techniques are implemented lying on blood donor's data set in favor of predicting the blood donor's attitude, which have been composed from the blood bank hub (Ankit Bhardwaj et al, 2012).

**B. Kamala** offers an analysis of require data mining facilities within cloud computing all along among situations of training underlying on the incorporated method of data mining with cloud computing. The data mining within cloud allows association on the way to consolidate the organization of software plus data storing space with guarantee of well-organized, dependable and safe services in favor of their users. The accomplishment of data mining methods throughout cloud computing will agree to the customers to get back significant evidence since it almost incorporated data warehouse that decreases the expenses of substructure as well storing space. This moves toward to decrease the fences that stay minor enterprises as of profiting of the data mining apparatuses (B. Kamala, 2012).

**Zeba Qureshi et al** throws light on association rule mining in cloud surroundings and a variety of similar and distributed mining algorithms. The accomplishment of association rule inside the distributed systems can be competently complete under Hadoop. Additional, the data move between the nodes and situations like node breakdown are concerned through Hadoop. This adds toughness and scalability in the direction of the system. A variety of parameters have an effect on the presentation of algorithms such as because time required to produce common item sets, bury the site communication cost, number of scans from side to side the database (Zeba Qureshi et al, 2011).

**Jayshree Jha et al** surveys the most relevant studies carried out in EDM using Apriori algorithm. Based going on the Apriori algorithm study and explore, this paper points away the major harms on the submission Apriori algorithm in EDM and presents an improved support-matrix based Apriori algorithm. The improved Apriori algorithm proposed in this research uses bottom up approach along with standard deviation functional model to mine frequent educational data pattern (Jayshree Jha et al, 2011).

**Zhang Chun-sheng et al** shows Two kinds of modified algorithm has been given based on the defect of the classical Apriori algorithm that cannot mine local credible rules: The first is Apriori-con algorithm based on confidence, this method proposed by the authors according to, the correctness and the application effect of the theory need to discuss by colleagues. Second is Apriori algorithm based on classification, the algorithm is further divided into Apriori-class-int algorithm based on interest classification, Apriori-class-int algorithm based on forecast classification, sorted Apriori-class-int algorithm based on clustering classification(Zhang Chun-sheng et al, 2011).

**Girish Kumar; Dr. Vibhakar Pathak** planned an original technique for optimization of rule mining using association rules. The proposed algorithm is a mixture of distance purpose and ant colony optimization. The paper describes a relation among a nearby big and worldwide big patterns that are used for local cutting at every site to minimize the examined candidates. Ant colony algorithm has been used for optimization of rule set(Girish Kumar; Dr. Vibhakar Pathak, 2011)

**Abhang Swati Ashok et al** shows that Apriori algorithm is used for finding best association rules. Apriori algorithm is aimed to control on transactional databases . Association rule mining is common as, assumed a group of itemsets, the Apriori algorithm efforts to discover subgroups which are frequent to smallest and least number of candidate (C) of the item sets. Apriori algorithm uses a "foot up" methodology. In which common subgroups are stretched one item at a period, and sets of candidates are verified in contradiction of the data. This algorithm ends when no other effective delays are established. The main purpose behind the Apriori Algorithm is to discover associations between dissimilar sets of data. Sometimes this is referred as Market Basket Analysis. Every set have multiple of items which called a transaction. The outcome of Apriori is group of rules that precise us that how often itemsets are enclosed in groups of data(Abhang Swati Ashok et al, 2011).

**M.Karunya et al** proposed Most important phase in analysis is data presentation. Data which is collected and signified in the two forms of datasets namely vertical or horizontal. To find out the dataset which is aggregated from horizontal database used three major methods like as CASE, SPJ, PIVOT. These techniques make the preprocessing of data simple with these query tools. In this paper horizontal aggregation

analysis is the process which combines and preprocesses data using Apriori algorithm(M.Karunya et al, 2011).

**G.Loshma et al** shows one of the most worth discussing theory of the data mining like as association rule finding for knowing about the habits of buying things from the market with respect to the customer behavior. Through this association rule mining frequent patterns are found over the dissimilar rules or transactions. As earlier all association rules are found by the horizontal data format. For finding association rules mainly two algorithms are used to find out frequent patterns from last many years which are Apriori algorithm and FP growth algorithm.

Few years ago some researchers found out the new approach through which vertical rules are achieved. In the horizontal association rule mining the rules are created using scanning database again and again which leads the transactions items support more complex. After some time this process is replicated for finding new rule from the candidate itemsets which increase the time complexity and scanning purpose more difficult. In this paper for the vertical approach the bitmap by this method memory saved and finding the frequent pattern become easy. In this paper both approaches vertical or horizontal are followed using multithreaded area. Author drives an improved parallel multithreaded method. This improved algorithm decreases the memory usage and protect time as procedure by bitmap demo of all datasets as well as bitmap algorithms(G.Loshma et al, 2010).

**Mohammed J. Zaki et al** shows that there are numbers of vertical process approaches are planned in features of association rule mining algorithm. By which efficient and effective results are shown by horizontal data approach. The chief benefit of the vertical data approach is its support for quicker frequency manipulative with connection processes on transactions as well as pruning widespread spontaneously of unconnected data. The main fault with all this methods is when middle outcome of vertical transactions become additional great for memory space, this outcome effect the procedures stability. In this paper a vertical data format representation named Diffset is proposed. Which take the record of different transaction ids of the candidate table pattern from generation of frequent patterns. Diffsets in this paper are cut down with the memory size need the storage for intermediate outputs. Paper through light that last vertical approach methods which used to increase the enactment significantly. A new

algorithm using diffsets for maximal mining patterns is proposed. The comparisons between these algorithms on the basis of dense moreover with sparse databases shows that diffsets provide magnitude performance enhancement over other previous technologies(Mohammed J. Zaki et al, 2010).

**K.Vanitha and R.Santhi** give reviews about an implementation of hash based Apriori algorithm. Authors examine hypothetically moreover experimentally principal of data structure regarding their solution. This kind of type of structure of data are the major factors in the efficient way of implementation. A new algorithm namely hash based algorithm created for candidate set generation. In this proposed algorithm it decreases the magnitude form the previous methods by resolving the bottleneck problem of performance. This approach improve the way of Apriori algorithm is as distinguished point that generation of shorter candidate sets allows the effective shortening of transactional database size at last past stages of the iterations which leads to reducing the cost in order of other iterations appropriately(K.Vanitha and R.Santhi, 2010).

**Sean Chester et al** covers common patterns planned whose jobs are to discover out collection of itemsets that frequently happened composed in the transactional databases at the profound of the datasets that produces quicker in speed access than the present algorithms to proceed regulator of datasets. Consequently, improvements are compulsory. This research paper covers the conclusive algorithm for the trouble by which the total vertical group drastically progress to its individualities of enactment when supervision of very enormous quantity of data is recorded. In this approach the practices of the average of enormous dataset weblogs after the FIMI 2004 conference to unlikeness among their demonstration alongside many state-of-the-art performances and regulate collected comparable efficiency with low-grade memory custom at thresholds furthermore besides the ability to mine confidence thresholds as unattempt in composed mechanisms. Authors are similarly designate how this determination can be extended to manage additional remarkable results(Sean Chester et al, 2010).

**K. Geetha** shows that common item generation is main method in zone of the association rule mining algorithms. The Data mining is demarcated as a technique of creating common itemsets that gratify smallest support. Well-organized algorithms to coalmine common patterns are vital in data mining field. The Apriori algorithm was basically proposed for generating the common itemset groups, which have been

numerous approaches planned to progress in its routine. But these approaches do not deal with the time constraint. Conversely, maximum among them still accept its candidate set generation-and-test methodology. In addition, many approaches do not produce all common patterns, creating them insufficient to originate association rules. The enhanced Apriori algorithm in this paper proposed needs less time in contrast to Apriori algorithm. So the time is decreasing while dealing with new approach(K. Geetha, 2010).

Maria S. Perez shows that Association rules are real treasured in addition extraordinary plans in abundant data mining circumstances. Apriori algorithm is the single of the famous association rule generation algorithm in data mining atmosphere . Apriori algorithm connects with a storage system in training to admission input in addition to output the results. This paper determines how to improve this technique explaining the main storage system to this difficulty through the routine of suggestions and parallel structures(Maria S. Perez, 2010).

**Sabita Barik, Debahuti Mishra** throws light on the key point of microarray research is to categorize genetic influence that are recorded in dissimilar techniques with esteem to divergent organic positions of cell in humankind and cases. Consequently, technique of data examination basics to be sensibly calculated such as grouping, organization, forecast etc. Authors in this paper have scheduled an well-ordered common pattern built on grouping to discover the gene which find common patterns screening analogous phenotypes foremost to exact symptoms for definite disease. Maximum of the methods for discovering common patterns were constructed on Apriori algorithm, that produces candidate gene sets step by step. This dispensation origins iterative dataset scans also great computational overheads. This algorithm likewise agonizes from drawing the support thresholds and confidence threshold structure to a crunchy limit. The hybridized Fuzzy FP-growth method not individual outperforms the Apriori with admiration to computational prices, then it similarly takes attention of scalability matters as the numeral of genes and disorder growths(Sabita Barik, Debahuti Mishra, 2009).

**Sheng Chai** shows the proficiency of excavating association rules is an significant meadow of awareness innovation in databanks. The Apriori algorithm approach is a traditional procedure in association rule mining. This paper offered a better-quality Apriori algorithm to upturn the competence of producing rules of association. This

procedure accepts a novel technique to condense the terminated group of sub-itemsets throughout pruning the candidate itemsets, that can produce straight the set of common itemsets in addition to eradicate candidates consuming a subsection that is not common in the meantime. This algorithm can increase the possibility of finding information in skimming databank and decrease the prospective balance of itemsets. Association rule mining is one of chief subjects of data mining investigation, a current and highlighting tool which chiefly selects the relative of dissimilar items in the databank. How to produce common itemsets is the main and basic query of frequent itemset mining. This is an significant feature in cultivating mining procedure that shows the shrinkage itemsets candidate in demand to produce common itemsets powerfully.

In traditional Apriori algorithm, while candidate itemset production are produced, the procedure wants to check their happening frequencies. The employment with termination outcome in great frequency in demanding, so marvelous quantity of capitals will be consumed whether in time or in memory.

An improved algorithm is projected for mining the association rules in producing common k-itemsets. In its place of trying whether these candidates are common itemsets afterward producing novel candidates, this fresh procedure discovers common itemsets straightly and eliminates the subset that is not common, which established on the traditional Apriori algorithm(Sheng Chai1, 2009).

**Rui chang** suggests an innovative optimization procedure named APRIORI-IMPROVE centered on the inadequate of Apriori. APRIORI-IMPROVE procedure offering optimization on 2-itemset generation, connection looseness and so on. APRIORIIMPROVE practices hash assembly to produce L2, customs an well-organized horizontal facts demonstration and enhanced approach of packing to protect time and memory. The enactment training displays that APRIORI-IMPROVE is ample quicker than Apriori. Authors have obtained a new procedure to mine all common itemsets, called APRIORIIMPROVE procedure. This procedure delivers contract databank by skimming simply one time and does not produce candidate sets. It detected that it is related to Apriori algorithm and FP-Growth displays ample advanced enactment improvement on sparse as fine as dense datasets(Rui chang, 2009).

**Yanxi Liu** throws light on the quick expansion of networks along with information technology, the limitless information have salaried additionally more consideration by individuals. Although in the search of material with great speed, the examination with mining of the evidence and guidelines secreted deep in the data are similarly paid further importance. Data mining equipment is to establish and investigate the data, which can remove to determine knowledge starting the mass of information, so in what way to relate the data mining equipment into innovativeness standard management is the attention of this theme planned. In this research paper, merging with data mining concept that is Apriori algorithm with respect to association rule mining procedure is designated in feature, the procedure employment procedure is demonstrated, along with the betterquality approaches of the algorithm are conversed. While current computer equipment with database machinery have been established speedily, might support the stock and speedily recover the impressive scale databanks or data storerooms, but then these methods were upto individual to collect these "massive" data, and not to successfully establish and expenditure the information secreted them, which ultimately led to today's occurrence of "ironic data which lead to deprived information". The development of data mining equipment met individual's needs. The equipment convoluted in non-natural intelligence, instrument wisdom, arithmetical analysis with other skills, and it creates decision examination into an original stage. In this paper the association rule mining techniques like as Apriori algorithm which is normally recycled in data mining is chiefly conversed(Yanxi Liu, 2008).

**Wei Wei, Songnian Yu** shows that inside the extent of mining association rules, preceding procedures, examples, FP-Growth algorithm and Apriori algorithm which have been commonly established with high judgements correspondingly. Many of these procedures decompose the problematic of mining association rules into dual sub problems: discover common pattern plus produce the wanted rules. Consequently, such a disintegration approach cannot but transport interval problematic when the proportions of catalog is significant plus creates user unendurable in a method where the necessary response time is précised. To crack the difficulty, faster and deep awareness of FPGrowth procedure and recommend an operative procedure by making the FP-tree structure which is known as Association Rule Growth that can concurrently find out common itemsets plus AR in a huge catalog. It is examined in concept that procedure is

precised along with association rules are produced by the procedure widespread. The experimentations display that the production of the AR arrangement produced by ARGrowth is narrowly true with the intervened runtime, in its place of the rapid epidemic in FP-Growth that is planned technique. On other side it lays a different on how to produce mining common patterns plus association mining rules at the similar time. On the extra pointer, it is selective to rage the methods for the limitation plus interest collection of rules. In presentation, it likewise supports OLAP schemes plus more or less intelligence supervisory software's that request gradual response. Planned an well-organized association rule growth technique, association rule growth, for concurrently mining common patterns plus association rules are used through FP-tree. With Equating the FP- Growth technique and AR-Growth technique can produce association rules normally it can be appropriate in with those presentations that necessitate direct responses. Coming investigation effort is positioned in how to impulse approximately limits into AR-Growth method to produce association that operators are fascinated in plus even build the involved supervisory schemes(Wei Wei, Songnian Yu,2008).



## CHAPTER 3

### SCOPE OF THE STUDY

---

Cloud Computing and Data Mining are the growing developments in the present planet of information expertise. Well known term Data Mining is a method of finding facts starting by the rare data in addition. Cloud computing delivers secure and stretchy communications which delivers the whole thing as a facility. By integrating data mining as well as cloud computing delivers quickness in addition rapid admission to expertise. The consequence of such mixing be supposed to be powerful and capacitive as to will be talented to contract through the growing manufacture of data, the circumstances for the well-organized mining of big quantity of data from a variety of data warehouses by means of the goal of making cooperative material or the manufacture of fresh awareness. The rising skill to make big quantity of data brings potential to find out and use precious knowledge of data. Process of Data mining is a winning method to examine data from dissimilar approaches and receiving helpful evidence since the data. It can support in forecasting styles or standards, categorization of data, labeling of data and in the direction of discover correlations, patterns as of the dataset. Today, the rising occurrence of cloud computing goals at converting the conventional technique of computing, by given that together software presentations plus hardware possessions as a facility. Enterprise information technology transportation suffers numerous expenses fluctuating from hardware overheads plus software preservation overheads to the overheads of observing, running an Information Technology infrastructure. The new arrival of cloud computing proposed a number of physical prospects of dipping some of those expenses;

nevertheless, concepts delivered by cloud computing are often insufficient to give main price reserves through the IT arrangement life-cycle.

Cloud substructure can be efficiently used for comprehensive and challenging operations with data that are distinctive for procedures of data mining. This term is essential to have obtainable accessible data warehouses in addition to accessible computing property that are able to believe. The accessible warehousing along with computing capitals ability delivers the well-organized way of storage as well as examining the great quantities of data.

## **CHAPTER4**

### **OBJECTIVES OF STUDY**

---

- Study and analysis of cloud using Cloud Sim simulator.
- Collect and Preprocessing of the data.
- Creating a Virtual cloud environment and interfacing of three algorithms on cloud.
- Apply and analyze the performance of Apriori and improved Apriori algorithm on the dataset.
- Apriori Algorithm based on horizontal data.
- FP Growth Algorithm based on projected data.
- Improved Apriori algorithm based on vertical data provides the very fast support count and probably with decrease execution time.
- Comparison of all three algorithm.

These objective are followed by the proposed work for the implementation. First point here is study and analysis of cloud using CloudSim simulator. CloudSim working is totally depend on the examples of CloudSim through which whole working of the CloudSim get clear. CloudSim is popular for the efficient working of it through the main 8 examples. First of all for the clear knowledge of the CloudSim have to clearly understand about the all examples. In all examples there are predefined attributes which work individually with the their respective working method. For instance in one example login procedure is followed in other registration procedure and so on. Collecting the dataset is very important in any proposed system as without dataset no implementation can be done completely. It is really important that the dataset cloud be

noise free and free from other outliers. Selecting the dataset and preprocessing it become necessary throughout the whole process. Preprocessing can be in many ways which includes filters and so on. One of the most important thing is to create virtual cloud environment through which work can be more efficient and effective. As nowadays cloud computing and data mining both collaborate to create the cloud mining approach which is in rising trends. Cloud computing changes the all work into more efficient way and enhanced way. As we know cloud provide us storage space . As in the proposed work the cloud approach is followed by the CloudSim working. Storage space is created on the local machine. By cloud usage user never to worry about storage space as all data saved on cloud storage area. After that algorithms of association rule mining namely Apriori algorithm and FP growth are surveyed. Apriori algorithm was the first algorithm developed for association rule mining. As all algorithm have parallel disadvantages over the advantages. As Apriori algorithm is known for the best frequent rule finding algorithm which leads to best algorithm. But the candidate generation of more tables leads to it disadvantages because it takes more time for scanning the database again and again for the generation of the candidate tables. Which overcome by the second algorithm named FP growth which removes the generation of candidate tables. It create the FP tree like structure through which all rules are uniquely identified. But this algorithm have also one disadvantage that maintaining the FP tree is more expensive. Which also lead to the more memory consumption. These two disadvantages are covered with the enhanced Apriori algorithm in which vertical approach of data mining Is followed which enhanced the results as well. All three algorithms after that implemented using CloudSim and outcomes of that are compared according to memory and time respectively.

Research for the implementation follows the following steps through which it shows all the implementation process take place. From the flowchart all main points regarding the implementation are well defined. By using this the process can be accurate and efficient rather than complicate.

### 5.1 Flowchart of work

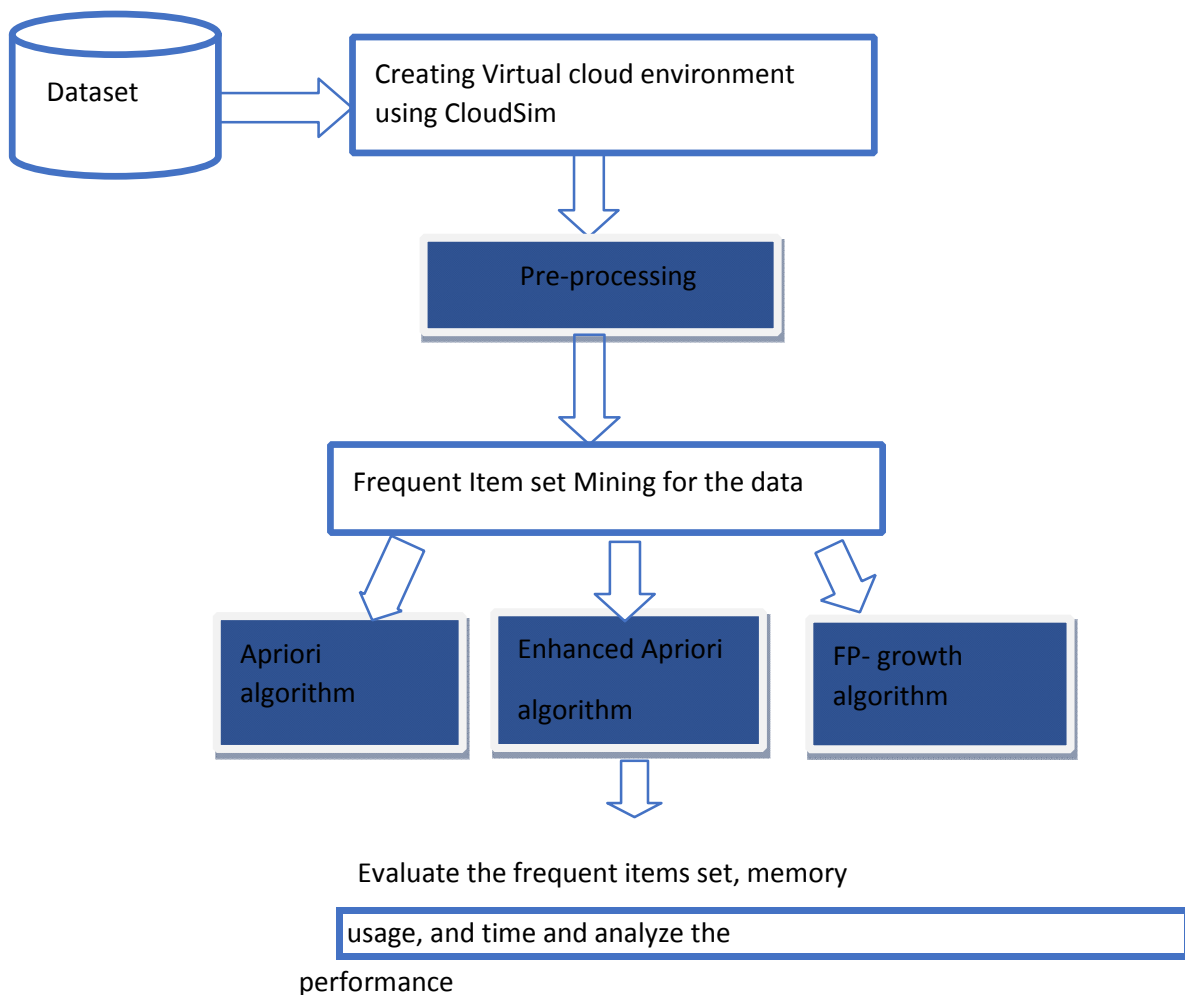


Figure 5.1:Flowchart of implementation

According to the flowchart first of all dataset is the first requirement for the implementation. Various dataset can be required for the checking purpose of

implementation. Only one dataset can also be good enough for the implementation it vary from kind of implementation. As my topic of thesis is association rule mining I use dataset of market analysis, supermarket, account etc. after that second step of the flow is creating virtual environment for the cloud. For creating it implementation use CloudSim in JAVA. Through CloudSim examples it's become easy to run and access .As for the virtual cloud we need some storage space we create that space in my local machine. Whenever we upload any dataset that on the cloud it firstly automatically saved to my local machine saved. And after that we can fetch it from the same storage space. All the process regarding saving the dataset become true by the CloudSim. After this step data preprocessing and collection is going to done. In which dataset selection done according to the requirement and missing value/noise removed using some preprocessing method. After that for finding the frequent itemsets three association rule mining algorithms are introduced through which frequent patterns are found. First algorithm is basic Apriori algorithm which uses the horizontal dataset format for finding the frequent itemsets. This algorithm is very basic algorithm for association rule mining which is in use from the last many decades. This algorithm is popular for its speed and accuracy. But it consumes lot of time in contrast for generating a lot of candidate generation tables. This difficulty is overcome by FP growth algorithm which neglect the generation of candidate tables. Through FP growth algorithm tree like structure is constructed. FP growth also take less time than the Apriori algorithm. But along with that in FP algorithm there was also one dis advantage like maintains of tree like structure became very costly and difficult.

## **5.2 Tool Used**

The proposed methodology is needed to be implemented in a simulation environment. The system is implemented in JAVA using Netbeans7 platform. CloudSim is printed in JAVA. The simply necessary information you require to make use of CloudSim is essential JAVA programming and a number of fundamentals concerning Cloud computing. Facts on programming IDEs such seeing that Eclipse or NetBeans is too cooperative. It is a library in addition to therefore CloudSim does not have on the road to be installed. Usually, you can unload the downloaded package in some directory, put in it to JAVA classpath and it is prepared to be used. Please confirm whether JAVA is obtainable on your system by organization JAVA –version and JAVAC authority.

### 5.3 Apriori algorithm

Apriori algorithm has been extremely commonly used for mining of common item sets and to find out associations. The main dissimilarity in Apriori algorithm is the less candidate itemsets it produces for analysis in every database go by. The exploration for association rule guidelines is showed by 2 key points: support plus confidence. Association rule mining Apriori results always return an association rule if and only if its support plus confidence morals are higher than user specified values. The outcome is organized through confidence. If numerous rulebooks have the related confidence then they all are supported by support and are the best rules. Thus Apriori favoritism extra best confident rules and describes these rules as extra remarkable .Apriori algorithm also uses a penetrating method called breadth first search method. It calculate the support of itemsets. It has a candidate generation fragment which makes use of the downward closure property of the support count. There are basic two chief steps of the Apriori algorithm are namely the join and prune steps.

(a) The first step is join step which is basically used to build fresh candidate sets tables. A candidate set is a itemset or table which can be either frequent or infrequent with regard to the support threshold given. Highest level candidate itemsets ( $C_i$ ) are produced by joining last level occurring itemsets are ( $L_{i-1}$ ) with new one.

(b) The prune step used in checking out and to decrease out candidate item-sets whose last levels are not common. This is mainly founded on the anti-monotonic property from which outcome of each subset of a common item set is also common. Furthermore a candidate item set table which is collection of one or many uncommon item sets of a last level is pruned from the method of common itemset and association rule mining.

Consider a transactional database in Table1 and assume that minimum support count is 2 and continue simple example of Apriori algorithm.

Database:

TID	Items
100	1,3,4
200	2,3,5
300	1,3,2,5
400	2,5

Table 5.1: generation of candidate tables from database

Candidate table1

Itemset	Support
1	2
2	3
3	3
5	3

Candidate table2

Itemset	Support
{1,2}	1
{1,3}	2
{1,5}	1
{2,3}	2
{2,5}	3
{3,5}	2

Candidate table3

Itemset	Support
{2,3,5}	2

#### 5.4 FP growth algorithm

Apriori algorithm always needs(  $i+1$ ) scans, where  $i$  is the length of the longest outline, we can use FP growth technique to condense the amount of scans of the entire transactional database  $D$  to find the frequent itemsets using only two scans of database.

**Example:** Consider a minor transactional database having six items as exposed in Table 5.2, figure1 gives the chief execution procedure of FP algorithm. This is explained as follows:

Steps for cracked example:

1. Scan the database and discover items with frequency larger than or identical to a threshold.
2. Order the frequent items in declining order {Q5, R4, P3, S3, T3}
3. Build a tree which has single root
4. Scan transactional database again; for each sample:
  - a) Improve the items from the sample to the current tree, using only the
  - b) Common items (i.e. items discovered in step 1.)
  - c) Recurrence (a). until all models have been process.

Table 5.2: Database for FP growth example

TID	Itemset
1	P,Q,R,S,U
2	Q,R,T
3	P,Q,S
4	P,Q,U
5	R,S,U
6	T,U

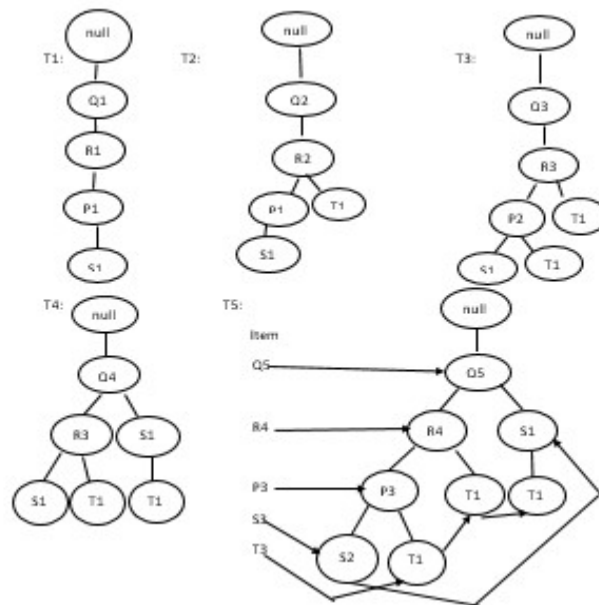


Figure 5.2:Fp tree

### 5.5 Improved Apriori algorithm

In this improved Apriori algorithm Vertical approach of association rule mining is used to find out common summaries from a straight transactional data organization which signifies the items categorized into specific transactions of itemsets ,in contrast data format of vertical method indicates as transactions categorized into specific items. Therefore, for a specific item, set of matching transaction IDs are present . As an alternative of straight data format, Apriori algorithm can be protracted to practice vertical approach data format for well-organized mining of association rule. Table 5.3 determine the data signified in Table 5.4, In the vertical data format. So that data can



be understood by Table 5.3, the vertical association rule mining approach design is the only one which analyze the dataset in one time to catch the set of itemsets. For the first itemset table generation afterward the second itemset generation, it simply needs to refer the earlier itemset . This rejects the circumstance to examine over the dataset each time to sum the happening of itemsets for all curved. Moreover, comparison with these algorithms well-known for the usage on transactional databases along with the horizontal data format of dataset, these algorithms which are used for the vertical approach tend to be extra finest. Some points which are discussing when the vertical data outline developed:

- Scanning the database and transform the dataset from horizontal format to the vertical format.
- For calculating the 1 candidate frequent item sets, the vertical t-id list data base is given and for each item it simply reads it corresponding trans id list from the given data base and incrementing the items support for each entry.
- Calculate the support count of an item set which is the length of id\_set of the item set.
- The frequent k-item sets can be used to make the candidate (k+1)-item sets using the Apriori based property.
- This procedure repeats, through k incremented by 1 every time, awaiting no common items or rejection candidate item sets can be establish.

Table 5.3: Database of horizontal format

Table 5.4: Conversion from horizontal to vertical database

TID	List of Items_id
D10	T1,T2,T3
D20	T1,T2,T3,T4,T5
D30	T1,T3,T4
D40	T1,T3,T4,T5
D50	T1,T2,T3,T4
D60	T2,T3,T4,T5
D70	T2,T3,T4
D80	T2,T4,T5
D90	T2,T4
D100	T2,T3
D110	T3,T4,T5
D120	T3,T4
D130	T3,T5
D140	T3,T4,T5
D150	T2,T3,T4,T5

ITEMS_ID	List of Tid	Support
T1	D10,D20,D30,D40,D50	5
T2	D10,D20,D50,D60,D70,D80, D90,D100,D150	9
T3	D10,D20,D30,D40,D50,D60, D70,D100,D110,D120,D140,D 150	13
T4	D20,D30,D40,D50,D60,D70,D 80, D90,D110,D120,D140,D150	12
T5	D20,D40,D60,D80,D110,D130, D140,D150	8

However the least support listed by the customer is 5 then for all the 1-itemset generations are common and while seeing table 1, In order to conclude the 2-candidate itemset is prominent then pruning is constructed on the support proposed. The 2-frequent itemset used for scanning the transactional databases over the many time which declines the amount of database scans necessary in order to compute the common item sets. This procedure is persistent until the number of common n-item-sets altered to zero. Optimization decays the many time scans of database which also decays the amount of time essential to compute the common item sets and respectively drops the space necessary to discover the common item sets and the efficiency of the vertical approach algorithm. For example in the Table 5.3 there are 12 items are present. In which Tid's with corresponding list of itemsets are present. Which both columns shows the database as a transactional table. In Table 5.4 the list of items changes to the vertical format in which first column occupy the list and after that in second column list of Tid's corresponding to that item sets. Additionally it includes the third column in which support is considered according to the occurrence of the itemsets. Through this method of changing the transactions into vertical format from the horizontal it enhances the results according to the time and memory prospective. And it decrease or completely neglect the database scanning again and again which leads to the improved results.

### 6.1 Experimental work

In direction to estimate the enactment of our projected algorithm, we have piloted experimentations on a PC (CPU: Intel(R) Core i3, 3.16GHz) with 16GByte of memory running Windows8.

In the implementation first of all registration form is created through which user can register by giving userID, username, password and also user role can be assigned according to the role what user want to do. After that clicking on the register button a user will be created.

#### Register user

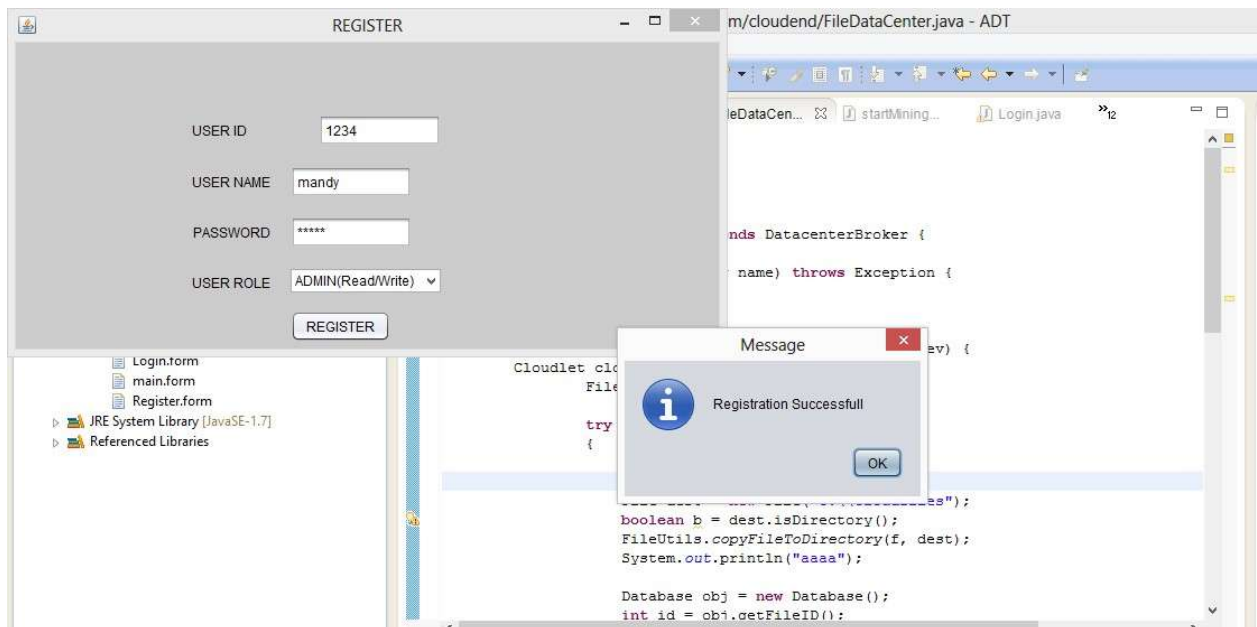


Figure 6.1: Register user page

In the above screenshot of register page, user id 1234 is assigned along with “mandy” as username and again “mandy” as password. Password is secured using (\*) rather than displaying “mandy” it displayed (\*\*\*\*\*). After that user role can be assigned like user can be admin with both (read/write)

And can be simple user with only read role. After clicking on the register button popup message of registration successfully is appeared . If the given information is not fully correct then the registration unsuccessful popup message will appear. All process of registration is going through CloudSim. This helps in enhancing the efficiency and time. By creating registration through CloudSim’s predefined attributes using main three files namely datacenter, data broker, virtual machines. Creating registration through these is enhancing the time and memory for the database. Mysql is use for the database saving of registration process. Table is created through which the registration process in connected.

### Login user

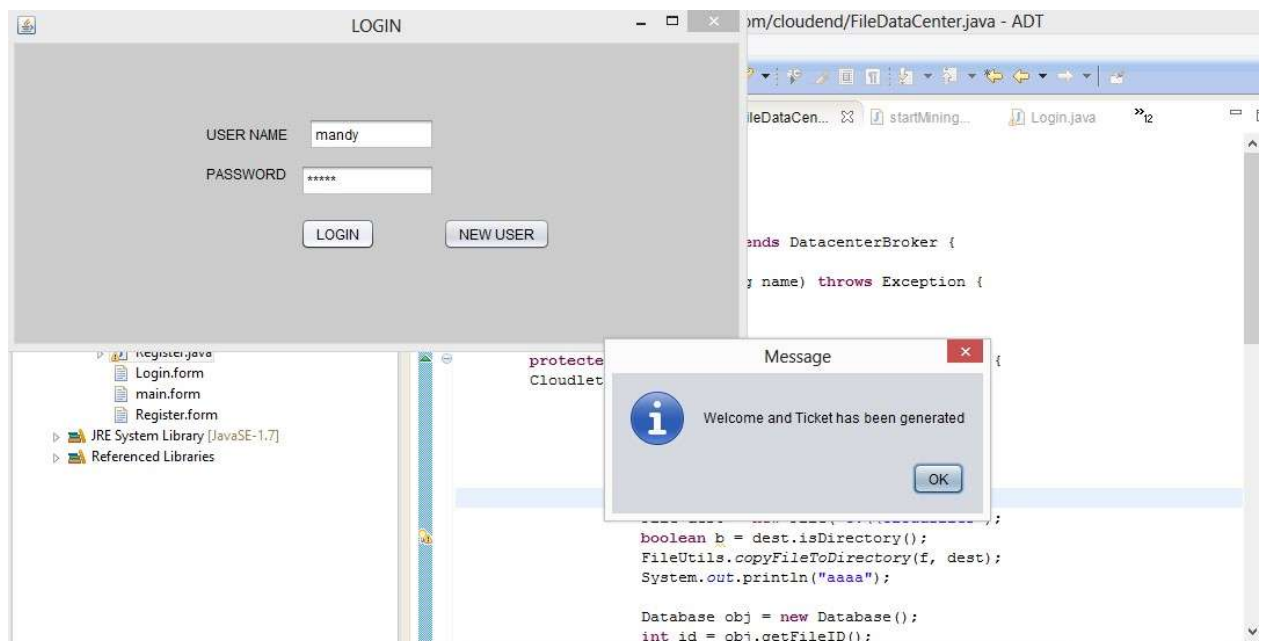


Figure 6.2: login user page

After registration process user can login with the same username and password which he/she given in the registration time. Like in previous screenshot username and password as “mandy”(\*\*\*\*\*) is given. User can login using that username and password into the login page. In the login page there is also one new user button which also have responsibility for creating new user through registration page as earlier process of registration is done. A message for successful login is

generated by displaying message “welcome and ticket has been generated”. If in the login page username and password given incorrect it leads to unsuccessful login message display. That wrong entry is not allowed for the login and lead to the unsuccessful login. Again login attributes are fetching through CloudSim’s files. Backend is connected with mysql in which table of login is created having columns username and password.

## Upload dataset and fetching

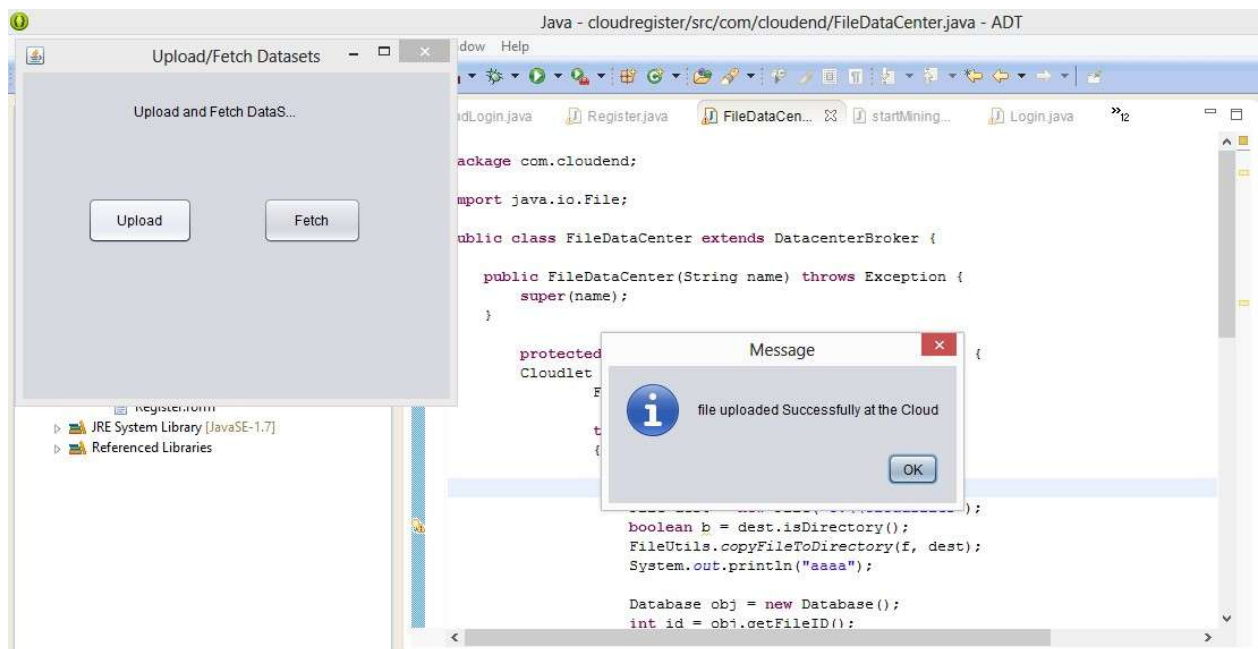


Figure 6.3: upload and fetch dataset

After login into the login page the upload/fetch datasets page is open as shown above. As the work is going through CloudSim the dataset uploading is done virtually on the cloud. Dataset is virtually uploaded on the cloud. As we created space into the C drive of our own system. When uploading of any dataset is done then it saved into the folder into the C drive of my machine. Once uploading of dataset is done on the cloud. For the instance my local C drive is my cloud space. After uploading the dataset onto the space. We can fetch it any time from the local cloud space. Once the upload process of dataset to the cloud is done message

regarding file uploaded successfully at cloud is popup. Again behind this CloudSim's three main file are running which are helping in the easy to make the code regarding upload process. At the backend mysql again upload process connected with a table in which file name, file size and type of file are saving.

### Selecting user for mining

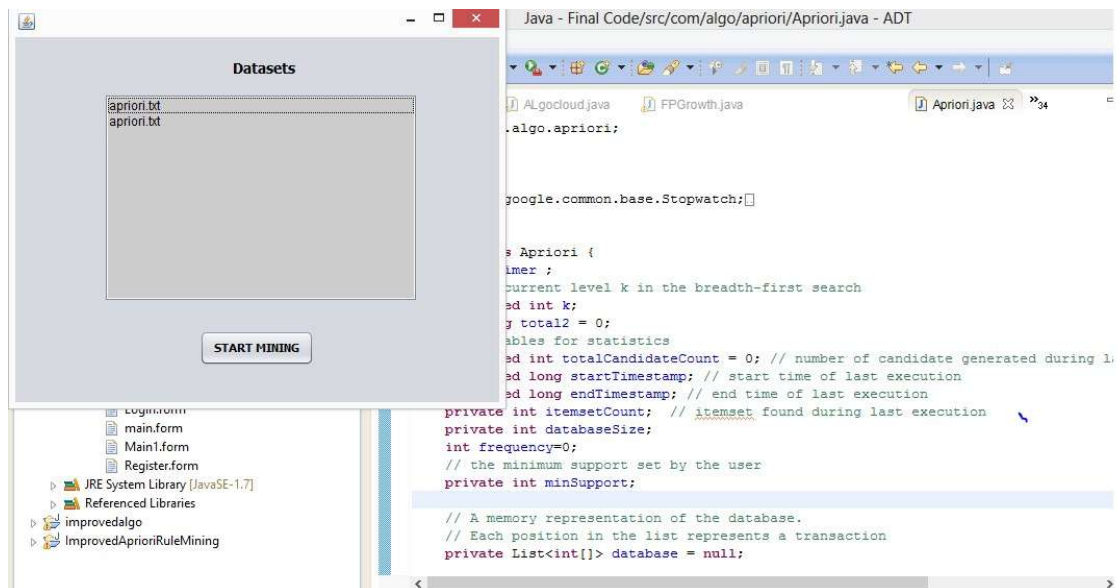


Figure 6.4: selection of user for data mining

As there are so many users created using registration here ,all users are displayed whether they can contain different user roles. After the uploading successful its necessary to choose one user who mine the dataset according to itself requirements. Here as above shown there is list of users is displayed. One of them can be choose for the mining process of association rules. As in above screenshot “mandy” as a user is selected for the mining. After the selection when we click to the start mining button. New main window will be opened. And also behind this the mysql queries are started running. And list of user is also fetched from the registration table. where all user are stored.

## 6.2 Data analysis and Interpretation

On the test dataset the below figure displays the outcomes of Horizontal Association Rule Mining Using Apriori algorithm, vertical association rule mining using Apriori algorithm and projected association rule mining using FP Growth:

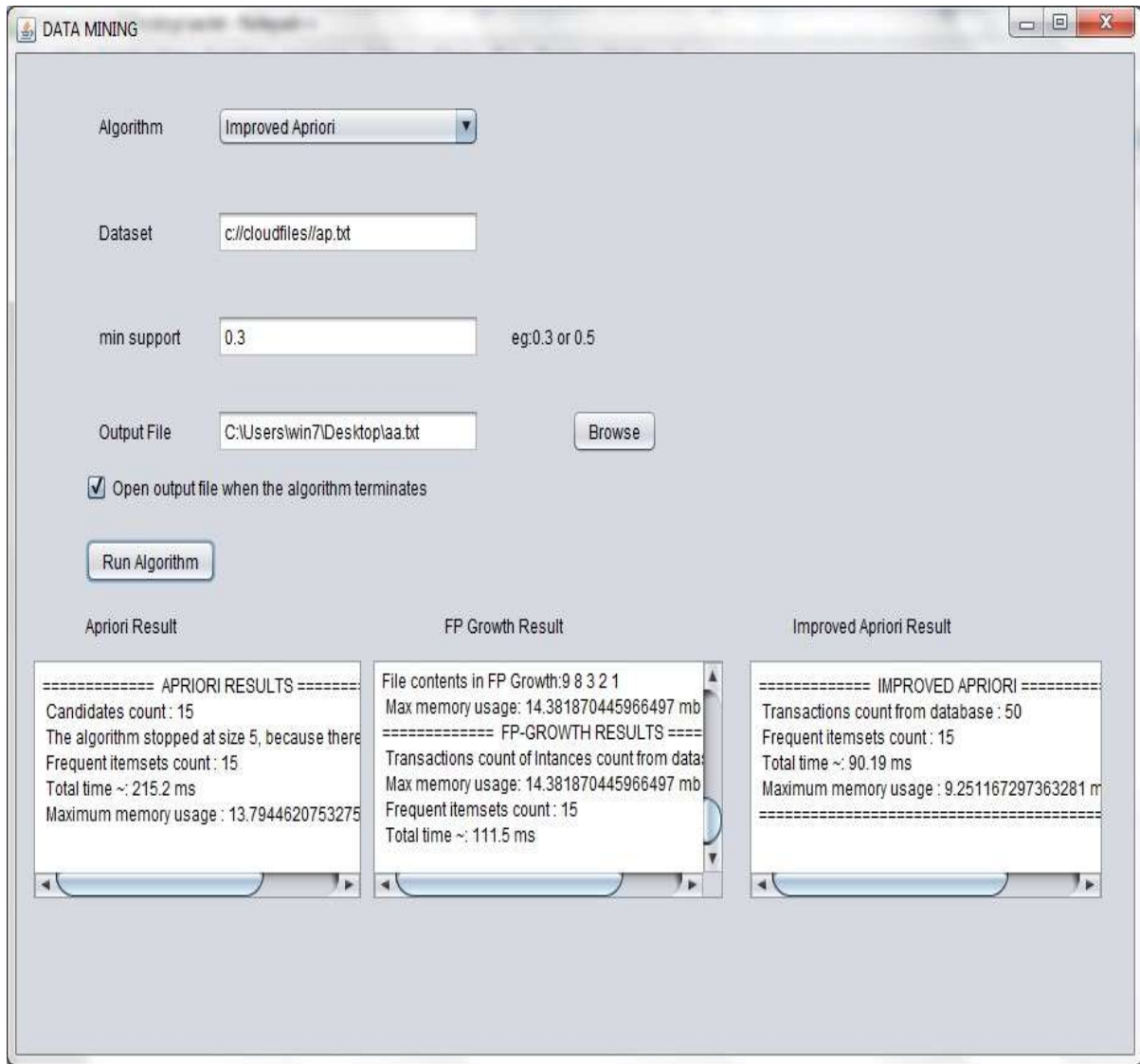


Figure 6.5: Results



### Performance Evaluation

The following bar graphs shows the Analysis of Time and memory with three association rule mining algorithms.

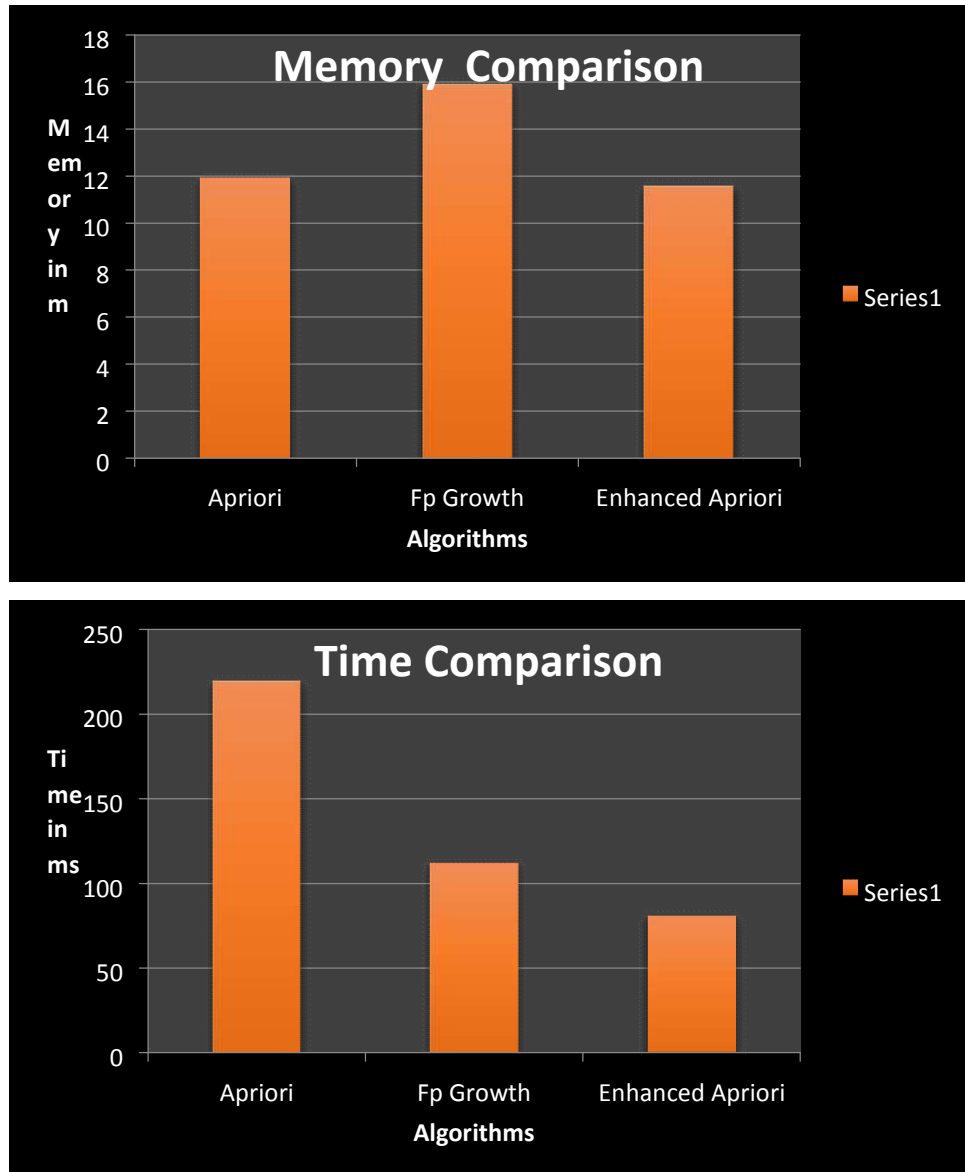


Figure 6.6: Analysis of Time and memory between Vertical Apriori ,Horizontal Apriori and FP growth algorithms

On the market basket dataset of market the below figure displays the outcomes of Horizontal Association Rule Mining Using Apriori algorithm, vertical association rule mining using Apriori algorithm and projected association rule mining using FP Growth:

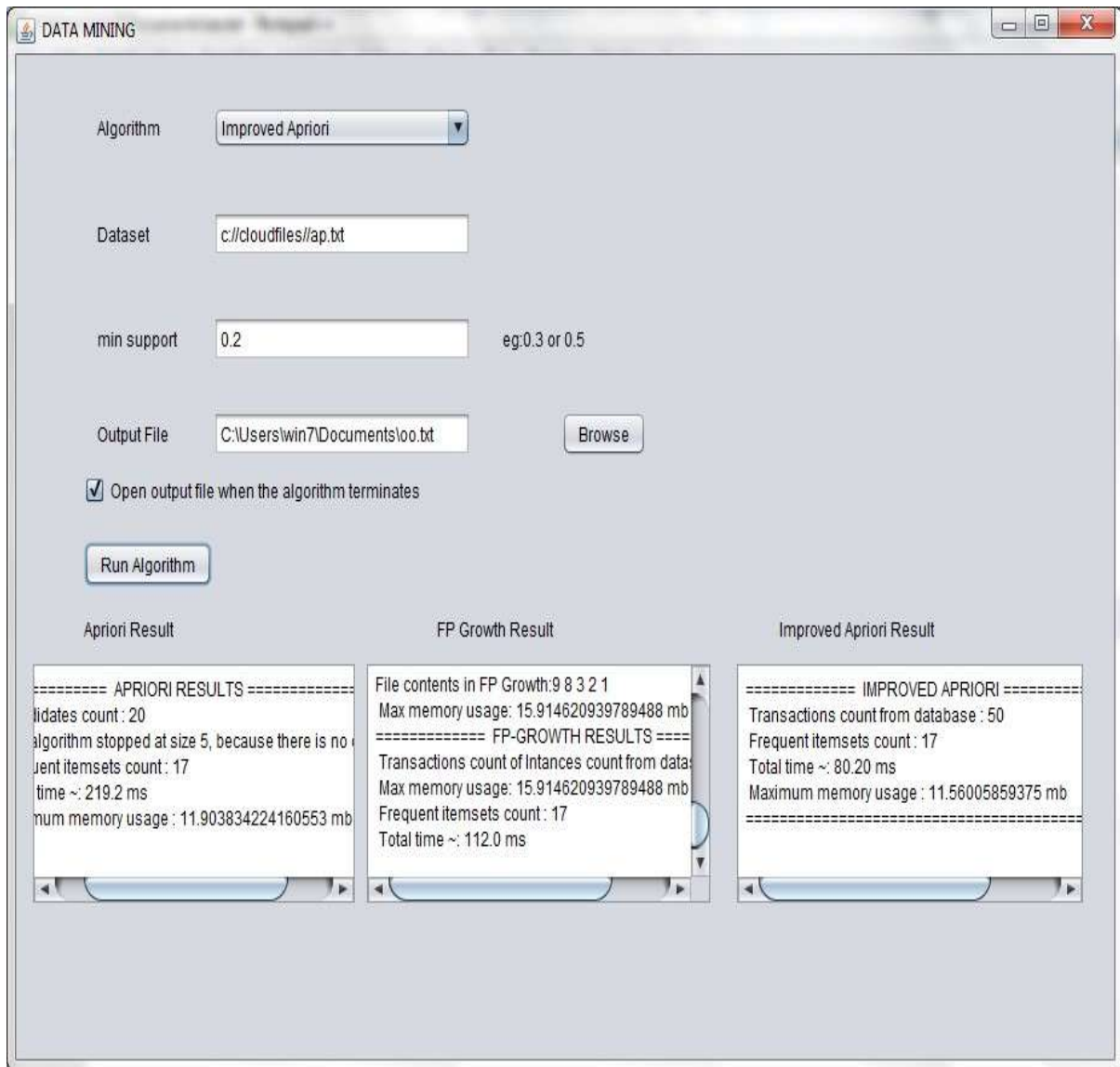


Figure 6.7: Results

### Performance Evaluation

The following bar graphs shows the Analysis of Time and memory with three association rule mining algorithms.

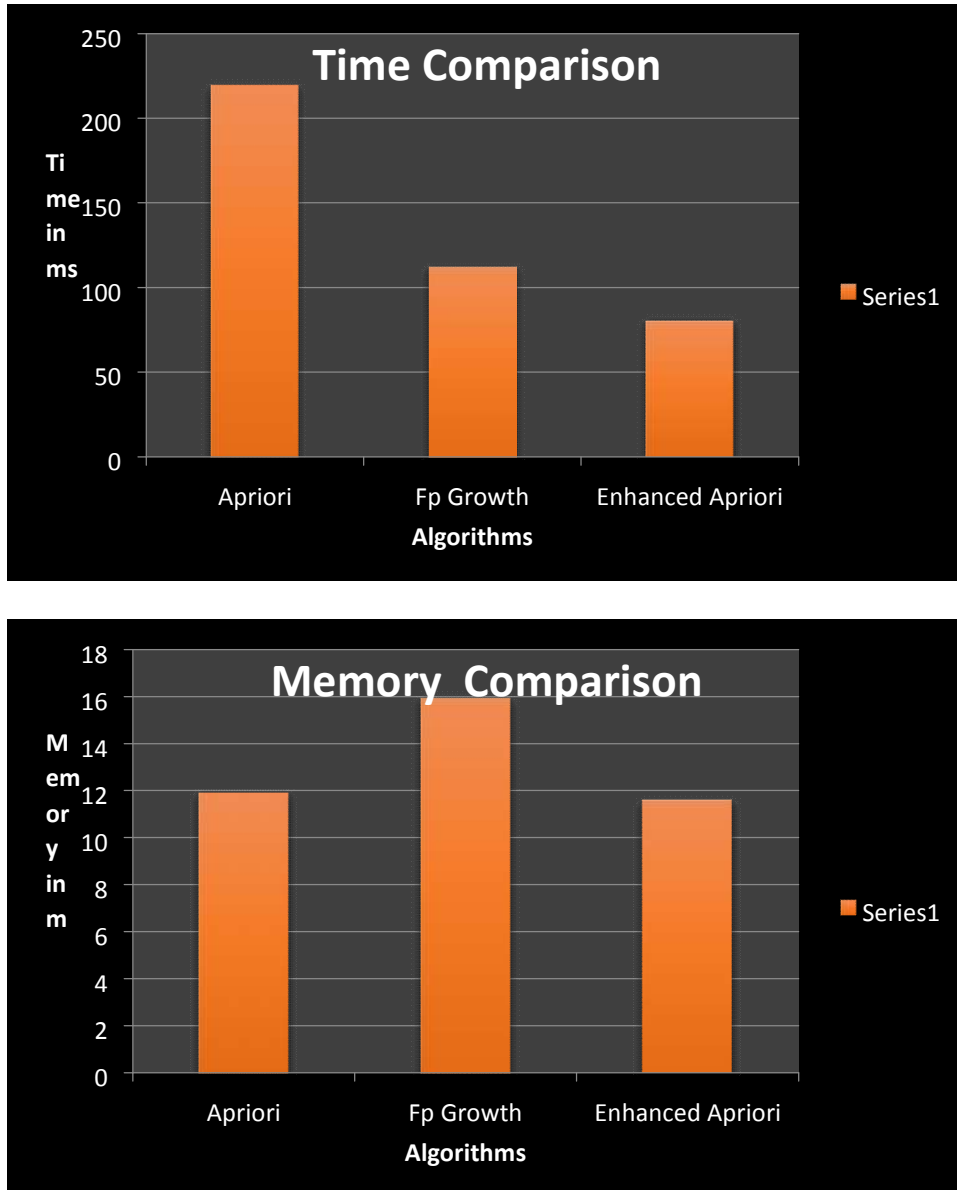


Figure 6.8: Analysis of Time and memory between Vertical Apriori ,Horizontal Apriori and FP growth algorithms

In the association rule mining region, the majority of the study efforts went into civilizing the algorithmic performance. Over the decade a mixture of algorithms that speak to these issues throughout the modification of look for strategies, pruning techniques as well as data structures have been developed. Whereas most algorithms' focal point on the open discovery of all policy that convince negligible support and confidence constraints for a specified dataset, growing consideration is being agreed to particular algorithms that try to get better processing time. Apriori algorithm finds the entire the common itemsets by scanning the database time following time, in addition it consumes a lot of point in time and memory space when scanning the database through mass data, which becomes the blockage of Apriori algorithm. The algorithms based on generated as well as tested candidate itemsets hold major drawbacks such as the database have to be scanned several times to produce candidate sets. Many scans will raise the I/O load and is timeconsuming. The production of vast candidate sets and computation of their support will put away a lot of CPU time. Apriori algorithm only extracts data based on horizontal data approach. While the FP growth algorithm is based on projected data approach. Both the algorithms pose few challenges. FP-Growth algorithm recursively generates vast amounts of provisional pattern bases plus conditional FP-trees at what time the dataset is vast. In such a case, together the memory procedure and computational cost are luxurious, such that, the FPtree cannot get together the memory condition.

Taking all these shortcomings keen on consideration, an enhanced Apriori algorithm has been planned in this effort. The improved Apriori algorithm is based on vertical data approach and improves the performance in terms of memory and time consumption.

In this practice, three association rule mining algorithms are proposed namely Apriori algorithm, FP growth algorithm and improved Apriori algorithm. Apriori algorithm use horizontal data format, FP growth use projection data format and improved using vertical data format. Using vertical format huge improvement can be seen in favor of time and memory.

**Reference to Books**

Heiko bock. Definitive guide to netbean platform 7.

Hand D., Mannila H. and Smyth P., Principles of Data Mining, MIT Press, 2001.

Tan P.-N., Steinbach M. and Kumar V., Introduction to Data Mining, Addison Wesley, 2006.

Michael J. Kavis, architecting the cloud: design decisions for cloud computing models

**Reference to web pages** <http://clean-clouds.com/2014/01/19/CloudSim-eclipse/>

[http://en.wikipedia.org/wiki/Apriori\\_algorithm](http://en.wikipedia.org/wiki/Apriori_algorithm)

[http://en.wikibooks.org/wiki/Data\\_Mining\\_Algorithms\\_In\\_R/Frequent\\_Pattern\\_Mining/](http://en.wikibooks.org/wiki/Data_Mining_Algorithms_In_R/Frequent_Pattern_Mining/)

[The\\_FP-Growth\\_Algorithm https://code.google.com/p/CloudSim/wiki/FAQ](https://code.google.com/p/CloudSim/wiki/FAQ)

**Reference to papers**

Abhang Swati Ashok, JoreSandeep S, “The Apriori algorithm: Data Mining Approaches Is To Find Frequent Item Sets From A Transaction Dataset”, International Journal of Innovative Research in Technology & Science(IJIRTS).

Ankit Bhardwaj, Arvind Sharma, V.K. Shrivastava, “Data Mining Techniques and Their Implementation in Blood Bank Sector –A Review”, International Journal of Engineering Research and Applications (IJERA) ISSN: 2248-9622.

B. KAMALA, “A Study On Integrated Approach Of Data Mining And Cloud Mining”

Chien-Hua Wang, Sheng-Hsing Liu and Chin-Tzong Pang, “Mining Association Rules Uses Fuzzy WeightedFP-Growth”, SCIS-ISIS 2012, Kobe, Japan, November 20-24, 2012.

Girish Kumar; Dr. Vibhakar Pathak, “An Improved Association Rule Mining Approach Using Distance Weight and Ant Colony Algorithm”, International Journal of Innovative Research in Technology & Science(IJIRTS)International Journal of Advances In Computer Science and Cloud Computing, ISSN: 2321-4058 Volume- 1, Issue- 2, Nov2013.

K. Geetha, Sk. Mohiddin, “An Efficient Data Mining Technique for Generating Frequent Item sets” International Journal of Advanced Research in Computer Science and Software Engineering, 2008.

K.Ganeshkumar H.Vignesh Ramamoorthy, S.Sudha D.Suganya Devi, “An Encrypted Technique with Association Rule Mining in Cloud Environment”, National Conference on Advancement of Technologies – Information Systems & Computer Networks (ISCON – 2012).

K.Vanitha and R.Santhi, “Using hash based Apriori algorithm to reduce the candidate 2-itemsets for mining association rule” Journal of Global Research in Computer Science .

Lingjuan Li, Min Zhang, “The Strategy of Mining Association Rule Based onCloud Computing”, 2011 International Conference on Business Computing and Global Informatization.

M.Karunya R.Reena Devi , “Frequent Pattern Analysis in Horizontal Layout Using Apriori Algorithm”, International Journal of Advanced Research in Computer Science and Software Engineering .

M.RaviKanth, G.Loshma, “Parallel multithreaded Apriori algorithm for vertical association rule mining”, International Journal of Advanced Research in Computer and Communication Engineering.

Mar'ia S. P'erez, Ram'on A. Pons, " An Optimization of Apriori Algorithm through the Usage of Parallel I/O and Hints".

Min Chen, XueDong Gao, HuiFei Li, "An Efficient Parallel FP-Growth Algorithm", 2009 IEEE.

Mohammed J. Zaki and Karam Gouda, " Fast Vertical Mining Using Diffsets", 2010 IEEE.

Nikita Jain, Vishal Srivastava, "Data Mining Techniques: A Survey Paper", IJRET: International Journal of Research in Engineering and Technology eISSN: 2319-1163 .ISSN: 2321-7308.

Qihua Lan, Defu Zhang, Bo Wu, "A New Algorithm For Frequent Itemsets MiningBased On Apriori And FP-Tree", 2009 IEEE.

Rui Chang, Zhiyi Liu, "An Improved Apriori Algorithm", 2011 International Conference on Electronics and Optoelectronics (ICEOE 2011).

Sabita Barik, Debahuti Mishra, "Pattern Discovery using Fuzzy FP-growth Algorithm from Gene Expression Data".

Sean Chester, Ian Sandler, Alex Thomo, " Scalable APRIORI-Based Frequent Pattern Discovery".

T.V.Mahendra ,Deepika ,Keasava Rao, "Data Mining for High Performance Data Cloud using Association Rule Mining", Volume 2, Issue 1, January 2012 ISSN: 2277 128X International Journal of Advanced Research in Computer Science and Software Engineering.

Tipawan Silwattananusarn and Dr. KulthidaTuamsuk, "Data Mining and Its Applications for KnowledgeManagement : A Literature Review from 2007 to2012", International Journal of Data Mining & Knowledge Management Process (IJDMP) Vol.2, No.5, September 2012.

Zeba Qureshi, Jaya Bansal, Sanjay Bansal, "A Survey on Association Rule Mining in

Cloud Computing”, International Journal of Emerging Technology and Advanced Engineering (ISSN 2250-2459, ISO 9001:2008 Certified Journal, Volume 3, Issue 4, April 2013).

Zhang Chun-sheng, Li yan, “Extension of Local Association Rules Mining Algorithm Basedon Apriori Algorithm”, 2014 IEEE.



**Abbreviations**

KDD: knowledge discovery database

GUI: graphical user interface

IT: information technology

APFT: Apriori- f tree

EDM: Educational data mining

AR: Association rules

OLAP: online analysis programming