# A review on an approach to mine significant itemset using tree based method and A-priori algorithm considering multiple attribute

**[1]Khushbu B. Joshi, [2]Prof. Smit M. Thacker**

[1]Research Scholar, [2]Assistant Professor
Department of Computer Engineering
HJD Institute of Technical Education & Research
Kera, Bhuj - Kutch

*Abstract*—Association rule mining or finding frequent patterns from massive set of associate degree owed knowledge has became a lively space in the sector of data discovery and diverse algorithms are developed to the present end. Traditional A-priori algorithm is used to mine a frequent pattern on the basis of Support-confidence however it doesn't consider about attributes like profit and quantity of an item. A modified approach makes use of traditional A-priori algorithm to generate a set of association rules from a database. Subsequently, the set of association rules mined are subjected to quantity (Q-factor) and profit (P-factor) to mine significant patterns. These factors are combined to induce PQ-gain on the basis of which association rules are mined. Another approach of segregating data by modifying the traditional A-priori algorithm is employed using tree based method. It helps to reduce the amount of space required to store the tables as well as time to mine frequent item set. The experimental results exhibit the adequacy and effectiveness of the modified approach in generating high utility association rules which can need fewer computations.

*Keywords*— *Data Mining , KDD , Association Rule , Market basket analysis, significant Item set, Support, Confidence, A-priori Algorithm, P- Factor , Q- Factor , PQ-Gain , Tree approach*

## I. INTRODUCTION

In today's emerging world, the role of data mining is increased day by day with the new aspect of business. Data mining has been proved as a very basic tool in knowledge discovery and decision making process. Data mining technologies are very frequently used in a variety of applications. Frequent item sets play an essential role in many data mining tasks that try to find interesting patterns from databases, such as association rules, correlations, sequences, episodes, classifiers, clusters. Data is important property for everyone as it disclosed the hidden information. This large amount of data from various sectors needs to be processed so that we can get useful information. Data mining is a technique to process data, select it, integrate it and retrieve some useful information [1]. Data mining is an analytical tool which allows users to analyze data, categories it and summaries the relationships among the data. It discovers the useful information from large amount of relational databases. Association rule mining discovers the interesting association or correlation among a large set of data items.

## II. ASSOCIATION RULE MINING

Association rules are if/then rules to reveal connections between apparently irrelevant information. Affiliation principle digging is learning for finding intriguing relations among the present variables from an extensive database. Association principle finds fascinating affiliation/connection among a substantial arrangement of information things. An illustration of association guideline would be "if a man is purchasing egg and bread together, he or she is liable to purchase margarine too". Today association rules are additionally utilized in numerous application ranges including web use mining, interruption identification.

### A. Concepts Of Association Rule Mining

*Item*: Field of a value-based database is known as items.

*Support*: It demonstrates how often times the things show up in the transaction (If there are X and Y items then count how many times they come together in transactions.)

$$Support(X, Y) = (XUY)/N \qquad [1]$$

N= Total number of transactions in database

*Confidence:* The certainty is the restrictive likelihood that, given X present in a move, Y will likewise be available.

$$Confidence (X, Y) = Support (XUY)/Support (X) \qquad [1]$$

*Frequent Itemset (Large Itemset (Li)):* The itemsets which fulfills the base support criteria are known as frequent itemsets.

Consider the following example:

TABLE I. TRANSACTIONAL DATABASE

| Transaction Id | Bread | Butter | Jam |
|---|---|---|---|
| 1 | 1 | 1 | 0 |
| 2 | 0 | 0 | 1 |
| 3 | 1 | 1 | 1 |
| 4 | 0 | 1 | 1 |

In Table I, 0 indicates the absence of an item and 1 indicates the presence of an item.

Considering X = Bread, Butter and Y = Jam

$$Support \{Bread, Butter\} => \{Jam\} = Support (X => Y)$$
$$= ¼$$
$$= 0.2 (20\%)$$

*Confidence {Bread, Butter}=>{Jam}* = Confidence( X=>Y)

$$= 0.2 / 2$$
$$= 0.1 (10\%)$$

## III. MARKET BASKET ANALYSIS

Market Basket Analysis is the best sample of affiliation tenet mining. This investigation recognizes the purchasing conduct of the client among different things that client places in their shopping cart [1]. The IDs of such client's conduct can help retailers to enhance the showcasing techniques to pick up the benefit into business. Market Basket examination is the method to determine relationship between datasets [1].

## IV. APRIORI ALGORITHM

Apriori, a standard ARM calculation is utilized as a part of the proposed way to deal with mine the association rules [2]. Apriori is one of the affiliation guideline mining calculation which is utilized to find all continuous itemsets from value-based database. This calculation utilizes earlier data of continuous thing set properties that is the reason it is known as Apriori calculation [1]. The procedure of affiliation principle mining utilizing Apriori is made out of two stages to be specific [2] :

- Frequent Itemset Generation: Generate all possible sets of relationship that have support value prominent than a predefined limit, called minimum support.
- Association Rule Generation: Generate association rules from the created successive itemsets that have certainty more noteworthy than a predefined limit called minimum confidence.

*APRIORI PROPERTY*: "All nonempty subsets of a frequent itemset must also be frequent" [3].

*A.Apriori Algorithm Steps*

*1)* First, the arrangement of competitor 1itemset is discovered (C1).
2) Then support is calculated by counting the occurrence of the item in transactional database.
3) After that we will prune the C1 utilizing least support Criteria. The thing which fulfills the base support criteria is thought about for the following procedure and which is known as L1.
4) Then again applicant set era is done and the 2-itemset which is produced known as C2.
5) Again we will figure the backing of the 2-Itemset (C2). What's more, we will prune C2 utilizing Minimum backing and produce L2.
6) This Process Continues till there is no Candidate set and continuous itemsets can be created.

## V. LITERATURE SURVEY

Arti Rathod et al.,[1] proposed a paper presents another methodology which remove noteworthy successive examples by applying so as to consider amount properties and Q-element and S-variable to the value-based database. Q-proportion is the proportion of amount of specific things all through all exchange to the aggregate amount of all things of all exchange and S-component is the result of Q-proportion of

specific things and the recurrence of specific things all through all exchange.

An improved approach takes Quantity attribute in the consideration. To mine significant (valuable) frequent patterns, count the Q-Ratio & S-factor of the association.

$$Q\text{-}Ratio = Q / \Sigma Q$$

Where, Q=Quantity of specific item throughout all transaction. ΣQ=Total Quantity of an all items presented in the transactional database.

$$S\text{-}factor = \Sigma Fi * Q\text{-}Ratio$$

S-factor is only a result of product of Q-ratio of particular association (items) and the recurrence of particular association (Items). Where, Fi=Total Number of occurrence or frequency of items or patterns in the Transactional database.

*A. Algorithm Steps of Improved approach:*

Consider the Transactional database where T=transactions, Tmax = Maximum available (total) transaction in database, M=size of Itemset in transaction T and maxM=maximum size of Itemset in transaction T.
1) Scan the Database and Initialize M=1, T=1.
2) Generate the association of transaction T of size M
3) Increment M by 1
4) If M<=maxM then GOTO step-2 else GOTO step-5
5) Increment T by 1
6) If T<=Tmax then initialize M=1 and GOTO step-2 else GOTO step-7
7) Combine all generated association in one table an eliminate the Duplicate generated association.
8) Count the Q-Ratio (equation (1)) of each generate association.
9) Count the Frequency of each generated association.
10) Count the S-factor (equation (2)) of each generated association.
11) Set all generated association in descending order.
12) END

The classical A-priori calculation has one disservice and that is it can just consider the vicinity and non attendance of the things and does not considers the importance of the things exhibited in the database such as the Quantity, weight and benefit characteristics. So the proposed work considers the Quantity quality from the value-based database and by applying Q-Ratio and S-variable equation we can discover the noteworthy successive examples. The enhanced methodology gives the huge continuous examples than the conventional Apriori calculation by considering Quantity characteristics.

Parvinder S. Sadhu and all., [2] proposed a productive methodology taking into account weight variable and utility for solid mining of huge affiliation rules. At first, the proposed approach makes utilization of the customary A-priori calculation to produce an arrangement of affiliation rules from a database. The proposed approach abuses the counter monotone property of the A-priori calculation, which expresses that for a k-itemset to be visit all (k-1) subsets of this itemset additionally must be visit. In this way, the arrangement of associations rules mined are subjected to weightage (W- gain) and utility (U- gain) requirements, and for each affiliation standard mined, a consolidated Utility Weighted Score (UW-Score) is registered. At last, a subset of

profitable affiliation rules in view of the UW-Score processed is being resolved.

The major steps involved in the proposed approach for association rule mining based on weightage and utility are:
Mining of association rules from D using -.
Step 1: Scan the Database.
Step 2: Computation of the measure W-gain.
Step 3: Computation of the measure U-gain.
Step 4:. Computation of UW-score from W-gain and U-gain.
Step 5: Determination of critical association rules taking into account  UW-score.

J.M. Lakshmi Mahesh [3] broke down the client conduct by applying so as to distinguish the recurrence set mix the systems and techniques. The proposed work use the specimen information gathered from the clients of a predefined Nationalized bank in India about its Retail benefits, &applying the A-priori idea for developing so as to recognize regular mining.The result has been represented as a frequency between various preferences, with Findings & Recommendations The predicted Findings and the hidden patterns/information can be applied for many purpose.

Raorane A.A. et al., [4]  analyzed the huge amount of data thereby exploiting the consumer behavior and make the correct decision leading to competitive edge over rivals. Market Basket Analysis is a device of learning revelation about co-event of ostensible or absolute things. Market Basket Transaction of Market Basket Analysis is an information mining procedure to determine relationship between information sets. Utilizing this proposed calculation the regular exchanges made by the clients have been broke down utilizing the backing and certainty of the clients in purchasing related things.

Lei Ji et al,.[6] taking into account the change on the established A-priori calculation, a high-measurement arranged A-priori calculation is proposed Unlike existed A-priori upgrades, our calculation receives another strategy to lessen the excess era of sub-itemsets amid pruning the competitor itemsets, which can get higher proficiency of mining than that of the first calculation when the measurement of information is high.

Huiying Wang et al., [7] calls attention to the bottleneck of traditional A-priori's calculation, displays an enhanced affiliation guideline mining calculation. The proposed calculation depends on diminishing the seasons of checking applicant sets and utilizing hash tree to store hopeful itemsets. The preparing time of mining is diminished and the effectiveness of calculation has progressed.

Ms Shweta et al., [8] utilizes A-priori to discover affiliation principle. The proposed work considers three affiliation guideline calculations: A-priori Association Rule, Predictive A-priori Association Rule and Tertius Association Rule. Author looks at the aftereffect of these three calculations and find that A-priori Association calculation performs superior to the Predictive A-priori Association Rule and Tertius Association Rule calculations.

Harveen Buttar et al., [9] proposed new calculation which is based upon the A-priori Algorithm that will improve the productivity and decrease time trait by making a model of model which will be valuable in beating the inadequacies of A-priori calculation. The proposed work checks a Database which will comprise of number of Items to be acquired by the client and aggregate Profit accomplished by the things .Profit proportion for every applying so as to thing will be figured Q-Factor .Profit Weight (PW-F) component is ascertained for those affiliation rules which are over the edge certainty.

$$Q\text{-}Factor = P / \sum P_i$$

Where i = 1 to n and $n$ = number of items and $P$ = profit of an item.

$$PW = \sum_{i=1}^{n} frequency * Q\text{-}Factor$$

The proposed work creates just number of interesting association designs that are both measurably and semantically critical for business improvement. Affiliation standard using so as to mine effectiveness can be enhanced traits like benefit, amount which will give the significant data to the client and also the business.

Abhijit Sarkar et al., [11] proposed another methodology of modifying so as to isolate information the customary A-priori calculation. The proposed approach minimizes the time many-sided quality, which gives great result.

### A.Tree Based Approach Algorithm

1. Support count of 1-itemset are checked and in view of the base bolster limit, a percentage of the itemsets are chosen. The chose itemsets are taken together to frame a set S, of length n, where n is the quantity of things.
2. The set S is considered as the root hub of the tree
3.The child nodes contain all concievable (n-1) length combinations, where n is the length of quick parent node.
4.Following the bottom-up approach, we figure the recurrence of the itemsets in the leaf nodes.
5. In the event that the recurrence of the itemsets in the leaf node is not exactly the predefined least support edge, then the leaf node and in addition its quick parent node is rejected.
6. In the wake of finishing the above steps, we get the last tree from which the chose nodes give the conceivable affiliation rules.

Reeti Trikha et al., [12] proposed a methodology by including new parameters which will be useful in giving most extreme benefit to the business associations. The proposed approach concentrates on the utility, centrality, amount and benefit of individual things for the mining of novel affiliation designs. The mined intriguing affiliation examples are utilized to offer profitable recommendations to an undertaking for strengthening its business utility.

### VI. CONCLUSION AND FUTURE WORK

In this proposed work we have seen how Association rule learning is a method for discovering interesting relations between variables in large databases. It is intended to identify strong rules discovered in databases using different measures of interestingness. A-priori algorithm is used to generate association rules from transactional database. Frequent Patterns

can be mined from association rules based on some interestingness measures.

A-priori doesn't consider the significance of other attributes like profit and quantity, so in future work different algorithms can be used to mine significant itemset based on multiple attribute and experimental results can be demonstrated to show the best algorithm.

## Acknowledgment

## References

[1] Arti Rathod "An Approach to Mine Significant Frequent Patterns by Quantity Attribute" 2014 IEEE Fourth International Conference on Communication Systems and Network

[2] Parvinder S. Sandhu, Atul Bisht, Dalvinder S. Dhaliwal, S. N. Panda," An Improvement in Apriori algorithm Using Profit And Quantity", Second International Conference on Computer and Network Technology,2010 IEEE

[3] J.M. Lakshmi Mahesh , " Association Models for market basket analysis, customer behavior analysis and business intelligence embedded with Apriori Concept ", International Journal of research in finance & Marketing, volume 2, Issue 1,january 2012, ISSN 2231-5985

[4] Raorane A. A , Kulkarni R.V and Jitkar B. D, "Associationr Rule- Extracting knowledge Using market basket analysis", Research Journal of Recent Sciences , Vo11 (2) 19-27, Feb. 2012.

[5] By Dr. M. Dhanabhakyam , Dr. M. Punithavalli " A Survey on Data Mining Algorithm for Market Basket", Global Journal of Computer Science and Technology, Volume 11 Issue 11 Version 1.0 July 2011.

[6] Lei Ji, Baowen Zhang, Jianhua Li, "A new improvement on apriori algorithm" Information Security Engineering School, Shanghai Jiaotong University, Shanghai, Chaina, 2006, IEEE.

[7] Huiying wang,Xiangwei Liu,"The Research of Improved Association Rules Mining Apriori Algorithm" 2011 IEEE Eighth International Conference on fuzzy Systems and Knowledge

[8] Ms Shweta, Dr. Kanwal Garg" Mining Efficient Association Rules Through Apriori Algorithm Using Attributes and Comparative Analysis of Various Association Rule Algorithms", Volume 3, Issue 6, June 2013 ISSN: 2277 128X International Journal of Advanced Research in Computer Science and Software Engineering.

[9] Harveen Buttar1, Rajneet kaur," Enhancement of Attributes of Apriori Algorit in Association Rule Learning", International Journal of Innovative Research in Computer and Communication Engineering Vol. 1, Issue 3, May 2013.

[10] Mamta Dhanda, Sonali Guglani, Gaurav Gupta," Mining Efficient Association Rules Through Apriori Algorithm Using Attributes", IJCST Vol. 2, Issue 3, September 2011.

[11] Abhijit Sarkar,Apurba Paul, Sainik Kumar Mahata, Deepak Kumar," Modified Apriori Algorithm to find out Association Rules using Tree based Approach",International Conference on Computing, Communication and Sensor Network (CCSN) 2012

[12] Reeti Trikha, Jasmeet Singh "Improvement in Apriori Algorithm with NewParameters ,International Journal of Science and Research (IJSR) 2012

[13] .Jiawei Han and Micheline Kamber "Data Mining : Concepts and Techniques" second Edition , ISBN 13 : 978-1-55860-901-3, ISBN 10: 1-55860-901-6, 2006 by Elsevier Inc