

Investigation and Application of Improved Association Rules Mining in Rapidminer

Zhi-hang Tang, Jian-hui Wei

School of Computer and Communications, Hunan Institute of Engineering, Xiangtan, Hunan, China
tang106261@126.com

Abstract — Association rule mining (ARM) is the process of finding frequent patterns and associations between set of objects from information repositories. Finding optimized techniques for generating association rules from large repositories has become a major area of study. The main purpose of this function is to find frequent patterns, associations and relationship between various database using different Algorithms. ARM is used to improve decisions making in the applications. ARM became essential in an information- and decision-overloaded world. They changed the way users make decisions, and helped their creators to increase revenue at the same time. Bringing ARM to a broader audience is essential in order to popularize them beyond the limits of scientific research and high technology entrepreneurship. ARM will assist you in reaching quality, informed decisions. Experimental results demonstrate that the improved FP-Growth algorithm is more efficient and suitable for the mining of large-scale transaction databases.

Keywords - Data Mining; Association rule mining; FP growth; knowledge discovery

I. INTRODUCTION

The practical insight about applications of association rule mining in various areas of day to day life. It provides the backend mechanism of data mining process, it shows the output of association rule towards the day to day life, and therefore it constructs a bridge between computer science and human life.

A typical and widely-used example of association rule mining is Market Basket Analysis. Data are collected using bar-code scanners in supermarkets. Such market basket databases consist of a large number of transaction records. Each record lists all items bought by some customers on a single purchase transaction. Managers of supermarket would be interested to know if certain groups of items are purchased together. They could use this data for adjusting store layouts, for cross-selling, for promotions, ARM will assist you in reaching quality, informed decisions.

II. RELATED WORK

Apriori algorithm has some limitation in spite of being very simple [1]. The major advantages of FP-Growth algorithm is that it uses compact data structure and eliminates repeated database scan FP-growth is faster than other association mining algorithms and is also faster than tree researching. According to [2] availability of determine "Which groups or sets of items are customer's likely to quality services is vital for the well-being of the economy. Market basket circles are covering all major aspects of the service analysis which may be performed on the retail data of customer transactions. This is the inherent cost of candidate generation, no matter what implementation technique is applied [4].

FP growth algorithm [5] generates frequent item sets from FP-Tree by traversing in bottom up fashion. It

allows frequent item set discovery without candidate item set generation. It is a two step approach.

Step 1: Build a compact data structure called the FP-tree .It is built using 2 passes over the data-set.

Step 2: Extracts frequent item sets directly from FP-tree. Traversal through FP-Tree Algorithm:

Input: A database DB, represented by FP-tree constructed and a minimum support threshold.

Output: The complete set of frequent patterns.

Method: call FP-growth (FP-tree, null).

Procedure FP-growth (Tree, a) {

1) **If Tree contains a single prefix path then // Mining single prefix-path FP-tree**

2) **Let P be the single prefix-path part of Tree;**

3) **Let Q be the multipath part with the top branching node replaced by a null root;**

4) **For each combination (denoted as β) of the nodes in the path P**

Do

5) **Generate pattern $\beta \cup a$ with support = minimum support of nodes in β ;**

6) **Let freq pattern set (P) be the set of patterns so generated;**

}

7) **Else let Q be Tree;**

8) **For each item a_i in Q do { // Mining multipath FP-tree**

9) **Generate pattern $\beta = a_i \cup a$ with support = a_i .support;**

10) **construct β 's conditional pattern-base and then β 's conditional FP-tree Tree β ;**

11) **If Tree $\beta \neq \emptyset$ then**

12) **Call FP-growth (Tree β , β);**

13) **Let freq pattern set (Q) be the set of patterns so generated;**

}

14) **Return (freq pattern set (P) \cup freq pattern set (Q) \cup (freq pattern set (P) \times freq pattern set (Q)))**

Advantages:

1) It uses Compact data structure.

2) It eliminates repeated database scan.

3) It is faster than Apriori algorithm.

4) It reduces the total number of candidate item sets by producing a compressed version of the database in terms of an FP tree.

Disadvantages:

- 1) It takes more time for recursive calls.
- 2) It is good only when user access paths are common.
- 3) It utilizes more memory

III. ASSOCIATION RULE MINING

One of the most common applications of ARM is market basket analysis [6] that discovers the relations among the items obtained by customers in the database. The improvement in the information technology allows all the retailers to obtain the daily transaction data at a very low cost [7]. Thus, the large amount of useful data to support the retail management can be extracted from large transactional databases. Data mining (DM) is used to obtain valuable information from large databases [8]. The aim of ARM analysis is to describe the most interesting patterns in an efficient manner [9]. ARM analysis (also known as the market basket analysis (MBA)) is method of determining customer obtained patterns by mining association from retailer transactional database [10]. Now a day's every product comes with the bar code. This data is rapidly documented by the business world as having the huge possible value in marketing. In detailed, commercial

organizations are interested in “association rules” that identify the patterns of purchases, such that the occurrence of one item in a basket will indicate the presence of one or more additional items. This “market basket analysis” result can then be used to recommend the combinations of the products for special promotions or sales, devise a more actual store layout, and give vision into brand loyalty and co-branding. It will also lead the managers towards efficient and real strategic decision making. Data mining (DM) methods are also used to find the collection of products, which are purchased together. It helps to choose which products should put side by side in the store shelves which may lead to important increase in sales. The problem of ARM can be decayed into the succeeding two stages [11].

A. Data Source

Figure 1, below, depicts a simplified relational model which might realistically be used by a supermarket to gather and store information about customers and the products they buy. It is simplified in that the attributes represented in each of the tables would likely be more numerous in an actual grocery store's database. However, to ensure that complexity of the related entities does not confound the explanation of Association Rules in this chapter, the tables have been simplified.

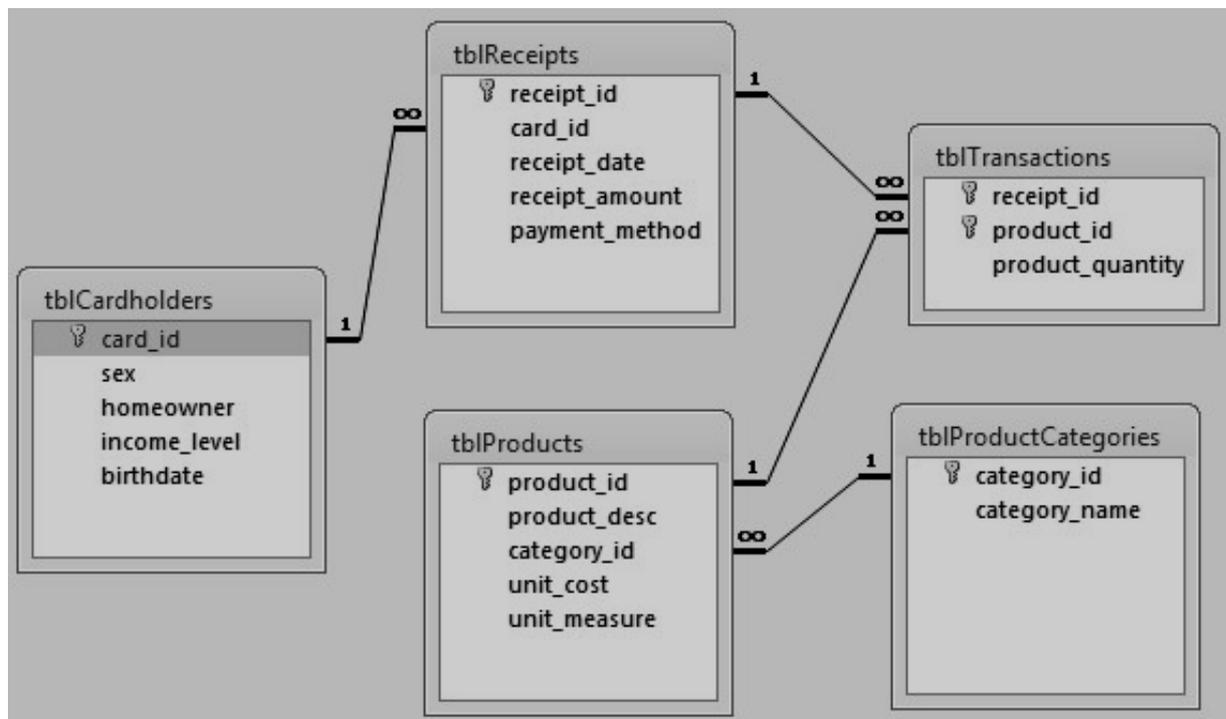


Fig.1 A simplified relational model of supermarket's database.

The datasets used throughout this paper consists of content and collaborative data. Content data was taken from the Supermarkets, Figure 2 depicts the first 19 rows of our previously discussed query, however this

query was run on tables containing 108,131 receipts from 10,001 different loyalty card holders, Figure 3 shows the Meta data view.

ExampleSet (108131 examples, 0 special attributes, 12 regular attributes)											View Filter (108131 / 108131): all		
Row No.	receipt_id	desserts	meats	juices	paper_goods	frozen_foods	snack_foods	canned_goo...	beer_wine_...	dairy	bread	produce	
1	1	0	1	1	0	1	0	0	0	0	0	1	
2	2	1	0	1	1	0	0	0	0	1	0	0	
3	3	1	1	1	1	1	0	1	1	1	1	1	
4	4	1	1	0	1	1	0	0	0	0	0	1	
5	5	0	0	0	0	0	1	0	1	0	0	0	
6	6	1	0	1	0	0	0	0	0	0	0	0	
7	7	1	0	0	1	0	0	0	0	0	0	0	
8	8	0	0	0	0	0	1	0	0	1	0	0	
9	9	1	0	0	1	0	0	0	0	0	0	0	
10	10	0	1	0	0	0	0	0	0	0	0	1	
11	11	0	0	0	0	1	0	0	0	0	0	0	
12	12	1	0	0	1	1	0	0	0	1	0	0	
13	13	1	0	0	0	0	0	0	0	0	1	0	
14	14	0	1	0	1	1	1	0	1	0	0	0	
15	15	1	0	0	1	0	0	0	0	0	0	0	
16	16	0	1	0	0	0	0	0	0	0	0	0	
17	17	0	0	1	0	1	0	0	1	1	0	0	
18	18	0	1	0	0	1	1	0	1	0	0	0	
19	19	1	0	0	0	1	0	1	1	0	0	1	

Fig.2 Query results from an expanded dataset

ExampleSet (108131 examples, 0 special attributes, 12 regular attributes)		
Role	Name	Type
regular	receipt_id	integer
regular	desserts	binominal
regular	meats	binominal
regular	juices	binominal
regular	paper_goods	binominal
regular	frozen_foods	binominal
regular	snack_foods	binominal
regular	canned_goods	binominal
regular	beer_wine_spirits	binominal
regular	dairy	binominal
regular	bread	binominal
regular	produce	binominal

Fig.3 the Meta data view

B. process of Association rule mining

Figure 4 depicts a basic operator workflow. Running the model on the entire dataset. If there are hundreds of thousands or millions of observations in your dataset, the model may take some time to run. Tuning the model

on a smaller sample can save time during development, and then once you are satisfied with your model, you can remove the sample operator and run the model on the entire dataset.

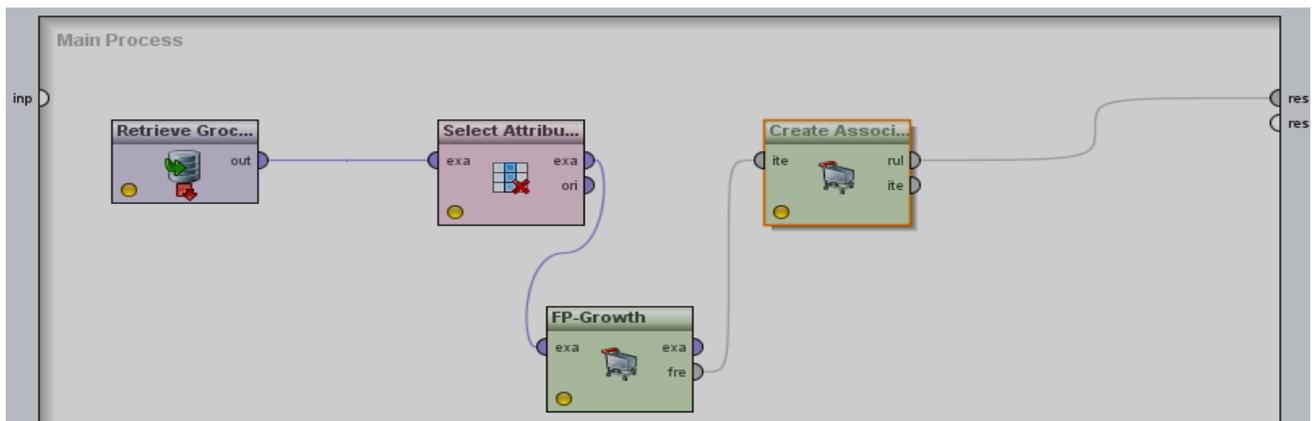


Fig.4 A basic Association rule mining operator workflow

Once any inconsistencies or other required transformations have been handled, we can move on to applying modeling operators to our data. The first modeling operator needed for association rules is FP-Growth (found in the Modeling folder). This operator, depicted in Figure 5, calculates the frequent item sets found in the data. Effectively, it goes through and

identifies the frequency of all possible combinations of products that were purchased. These might be pairs, triplets, or even larger combinations of items. The thresholds used to determine whether or not items are matches can be modified using the tools on the right-hand side of the screen.

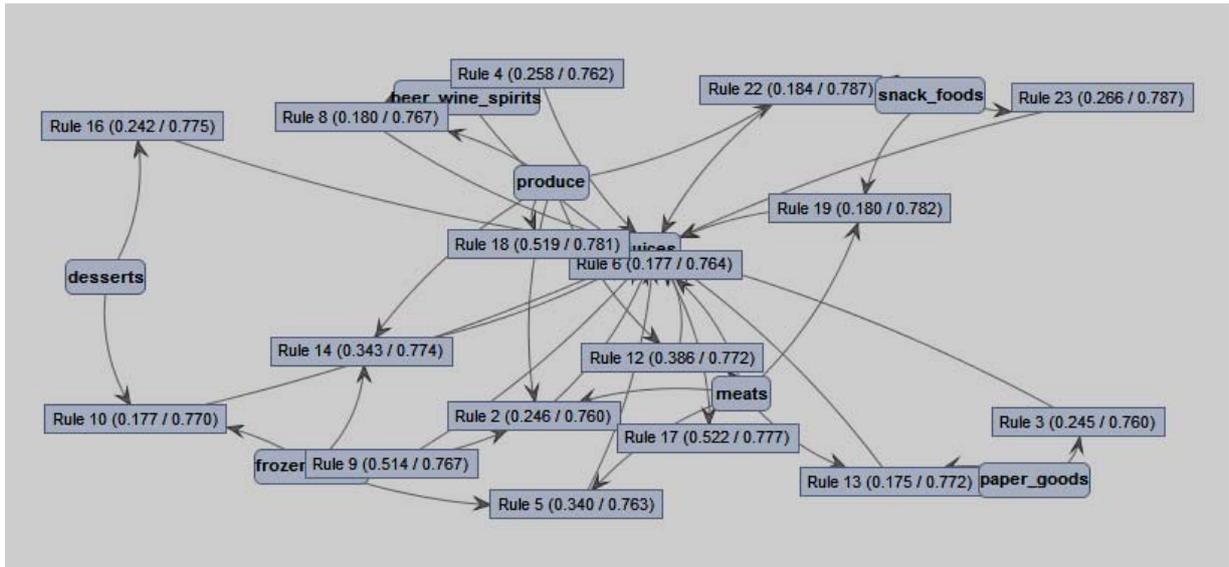


Fig.8 juices of Association Rules

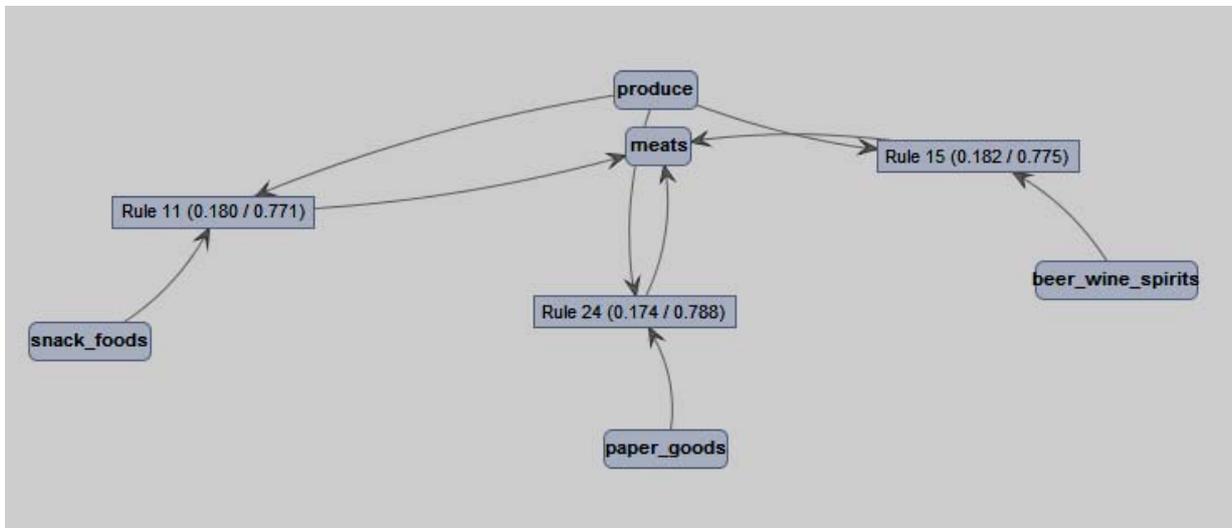


Fig.9 meats of Association Rules

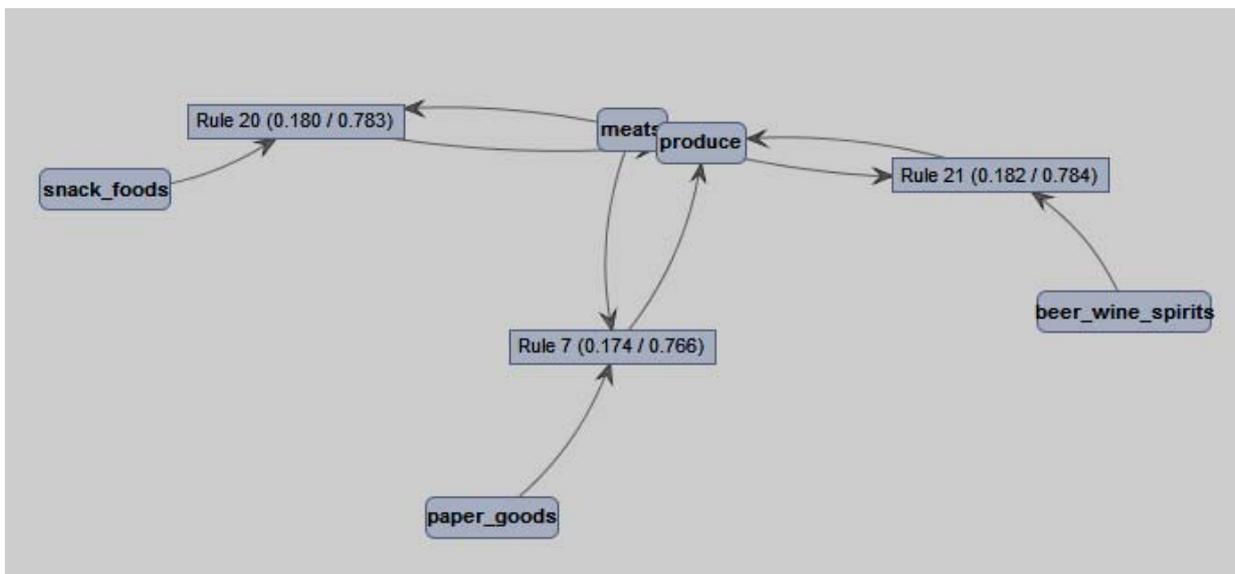


Fig.10 produce of Association Rules

When min support=0.5, Association Rules are as follows:

[meats] --> [juices, frozen_foods] (confidence: 0.506)
 [frozen_foods] --> [juices, meats] (confidence: 0.508)
 [frozen_foods] --> [juices, produce] (confidence: 0.512)
 [produce] --> [juices, frozen_foods] (confidence: 0.516)
 [beer_wine_spirits] --> [juices, meats] (confidence: 0.523)
 [beer_wine_spirits] --> [juices, produce] (confidence: 0.532)
 [snack_foods] --> [juices, meats] (confidence: 0.533)
 [snack_foods] --> [meats, produce] (confidence: 0.534)
 [beer_wine_spirits] --> [meats, produce] (confidence: 0.537)
 [paper_goods] --> [meats, produce] (confidence: 0.539)
 [paper_goods] --> [juices, meats] (confidence: 0.543)
 [snack_foods] --> [juices, produce] (confidence: 0.545)
 [meats, frozen_foods] --> [juices, produce] (confidence: 0.551)
 [frozen_foods, produce] --> [juices, meats] (confidence: 0.554)
 [beer_wine_spirits] --> [snack_foods] (confidence: 0.563)
 [snack_foods] --> [beer_wine_spirits] (confidence: 0.564)
 [desserts] --> [juices, frozen_foods] (confidence: 0.567)
 [meats] --> [juices, produce] (confidence: 0.574)
 [produce] --> [juices, meats] (confidence: 0.580)
 [juices, meats, produce] --> [frozen_foods] (confidence: 0.637)
 [meats, produce] --> [frozen_foods] (confidence: 0.646)
 [juices, meats] --> [frozen_foods] (confidence: 0.652)
 [juices] --> [frozen_foods] (confidence: 0.659)
 [juices, produce] --> [frozen_foods] (confidence: 0.661)
 [frozen_foods] --> [produce] (confidence: 0.661)
 [juices, frozen_foods] --> [meats] (confidence: 0.662)
 [meats] --> [frozen_foods] (confidence: 0.663)
 [frozen_foods] --> [meats] (confidence: 0.665)
 [juices] --> [produce] (confidence: 0.666)
 [produce] --> [frozen_foods] (confidence: 0.667)
 [desserts] --> [meats] (confidence: 0.667)
 [juices, frozen_foods] --> [produce] (confidence: 0.668)
 [juices] --> [meats] (confidence: 0.670)
 [juices, snack_foods] --> [meats] (confidence: 0.677)
 [desserts] --> [produce] (confidence: 0.680)
 [snack_foods] --> [meats] (confidence: 0.682)
 [paper_goods] --> [produce] (confidence: 0.685)
 [beer_wine_spirits] --> [meats] (confidence: 0.685)
 [juices, beer_wine_spirits] --> [meats] (confidence: 0.687)
 [juices, snack_foods] --> [produce] (confidence: 0.692)
 [snack_foods] --> [produce] (confidence: 0.692)
 [beer_wine_spirits] --> [produce] (confidence: 0.693)
 [juices, beer_wine_spirits] --> [produce] (confidence: 0.698)
 [paper_goods] --> [meats] (confidence: 0.704)
 [paper_goods] --> [frozen_foods] (confidence: 0.707)
 [juices, paper_goods] --> [meats] (confidence: 0.715)
 [juices, frozen_foods, produce] --> [meats] (confidence: 0.716)
 [juices, meats, frozen_foods] --> [produce] (confidence: 0.722)
 [meats, frozen_foods] --> [produce] (confidence: 0.725)
 [frozen_foods, produce] --> [meats] (confidence: 0.730)
 [juices, desserts] --> [frozen_foods] (confidence: 0.732)
 [desserts] --> [frozen_foods] (confidence: 0.737)
 [juices, meats] --> [produce] (confidence: 0.739)
 [juices, produce] --> [meats] (confidence: 0.744)
 [meats] --> [produce] (confidence: 0.744)
 [produce] --> [meats] (confidence: 0.752)
 [meats, frozen_foods, produce] --> [juices] (confidence: 0.760)
 [paper_goods] --> [juices] (confidence: 0.760)
 [beer_wine_spirits] --> [juices] (confidence: 0.762)
 [meats, frozen_foods] --> [juices] (confidence: 0.763)
 [meats, beer_wine_spirits] --> [juices] (confidence: 0.764)
 [meats, paper_goods] --> [produce] (confidence: 0.766)
 [produce, beer_wine_spirits] --> [juices] (confidence: 0.767)
 [frozen_foods] --> [juices] (confidence: 0.767)
 [frozen_foods, desserts] --> [juices] (confidence: 0.770)
 [produce, snack_foods] --> [meats] (confidence: 0.771)
 [meats, produce] --> [juices] (confidence: 0.772)
 [meats, paper_goods] --> [juices] (confidence: 0.772)
 [frozen_foods, produce] --> [juices] (confidence: 0.774)
 [produce, beer_wine_spirits] --> [meats] (confidence: 0.775)

[desserts] --> [juices] (confidence: 0.775)
 [meats] --> [juices] (confidence: 0.777)
 [produce] --> [juices] (confidence: 0.781)
 [meats, snack_foods] --> [juices] (confidence: 0.782)
 [meats, snack_foods] --> [produce] (confidence: 0.783)
 [meats, beer_wine_spirits] --> [produce] (confidence: 0.784)
 [produce, snack_foods] --> [juices] (confidence: 0.787)
 [snack_foods] --> [juices] (confidence: 0.787)
 [produce, paper_goods] --> [meats] (confidence: 0.788)

When min support=0.75, Association Rules are as follows:
 [produce] --> [meats] (confidence: 0.752)
 [meats, frozen_foods, produce] --> [juices] (confidence: 0.760)
 [paper_goods] --> [juices] (confidence: 0.760)
 [beer_wine_spirits] --> [juices] (confidence: 0.762)
 [meats, frozen_foods] --> [juices] (confidence: 0.763)
 [meats, beer_wine_spirits] --> [juices] (confidence: 0.764)
 [meats, paper_goods] --> [produce] (confidence: 0.766)
 [produce, beer_wine_spirits] --> [juices] (confidence: 0.767)
 [frozen_foods] --> [juices] (confidence: 0.767)
 [frozen_foods, desserts] --> [juices] (confidence: 0.770)
 [produce, snack_foods] --> [meats] (confidence: 0.771)
 [meats, produce] --> [juices] (confidence: 0.772)
 [meats, paper_goods] --> [juices] (confidence: 0.772)
 [frozen_foods, produce] --> [juices] (confidence: 0.774)
 [produce, beer_wine_spirits] --> [meats] (confidence: 0.775)
 [desserts] --> [juices] (confidence: 0.775)
 [meats] --> [juices] (confidence: 0.777)
 [produce] --> [juices] (confidence: 0.781)
 [meats, snack_foods] --> [juices] (confidence: 0.782)
 [meats, snack_foods] --> [produce] (confidence: 0.783)
 [meats, beer_wine_spirits] --> [produce] (confidence: 0.784)
 [produce, snack_foods] --> [juices] (confidence: 0.787)
 [snack_foods] --> [juices] (confidence: 0.787)
 [produce, paper_goods] --> [meats] (confidence: 0.788)

III. CONCLUSIONS

The study shows that it is not only valuable for business supervision but also helpful in medical and IT industry. This paper will convey strong source for implementation and use of association rule mining ARM to comfort the process of finding trends in large data sets and will be helpful in effective and efficient decision making. The main purpose is to find the frequent patterns (FP), association, and relationship between various databases using different methods. For example in market basket analysis(MBA) ,by utilizing the association rules(AR) that are generated as a result of analyses, the retail store manager will be able to expand and apply effective marketing strategies and in disease identification frequent patterns are generated to discover the frequently occur diseases in a definite area . The conclusion in all applications is some kind of association rules (AR) that are useful for efficient decision making.

ACKNOWLEDGMENT

This Project isSupported by Scientific Research Fund of Hunan Provincial Education Department (15A043).

REFERENCES

- [1] B. Santhosh Kumar and K.V. Rukmani. Implementation of Web Usage Mining Using APRIORI and FP Growth Algorithms, *Int. J. of Advanced Networking and Applications*, vol.6, pp. 400-404, 2010
- [2] Ankur Mehay, Dr. Kawaljeet Singh, and Dr. Neeraj Sharma. AnalyzeMarket Basket Data using FP-growth and Apriori Algorithm, *International Journal on Recent and Innovation Trends in Computing and Communication*, vol.16, pp. 693-696, 2013
- [3] Jiawei Han, Jian Pei, Yiwen Yin And Runying MAO. Mining Frequent Patterns without Candidate Generation: A Frequent-Itemset Mining Approach, *Data Mining and Knowledge Discovery*, vol.8, pp.53–87, 2004
- [4] Zhang H, Padmanabhan B, Tuzhilin A. On the discovery of significant statistical quantitative rules. In: *Proceedings of 10th international conference on knowledge discovery and data mining (KDD 2004)*, pp. 374–383, 2004
- [5] P. Yang and Z. Song. An improvement to FP-growth algorithm, *Journal of Anhui Institute of Mechanical & Electrical Engineering: Natural Science*, vol.17, pp. 8-13, 2005
- [6] Hahsler, M., Grün, B., & Hornik, K. Arules: a computational environment for mining association rules and frequent item sets. *Journal of Statistical Software*, vol.14, pp.1–25, 2005
- [7] Tan, P. N., Steinbach, M., & Kumar, V. *Introduction to data mining*. Reading: Addison-Wesley, 2005
- [8] Chen, M., & Lin, C.. A data mining approach to product assortment and shelf space allocation. *Expert Systems with Applications*, vol.32, pp.976-986, 2007
- [9] Tan, P., Steinbach. & Kumar, V. *Introduction to data mining*. Boston: Pearson Education, 2006
- [10] Tang, K., Chen, Y., & Hu, H.. Context-based market basket analysis in a multiple-store environment. *Decision Support Systems*, vol.45, pp.150-163, 2008
- [11] Agrawal, R., Imielinski, T., & Swami, A. Mining association rules between sets of items in large databases. *ACM SIGMOD Conference*, Washington DC, USA, 1993.