Comparative Study of Improved Association Rules Mining Based On Shopping System

Tang Zhi-hang

School of Computer and Communication, Hunan Institute of Engineering Xiangtan 411104, China

Email:tang106261@126.com

-----ABSTRACT-----

Data mining is a process of discovering fascinating designs, new instructions and information from large amount of sales facts in transactional and interpersonal catalogs. The main purpose of this function is to find frequent patterns, associations and relationship between various database using different Algorithms. Association rule mining (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. It 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.

Keywords: Comparative study, Association rule mining, FP growth, decision making

Date of Submission: Feb 02, 2016

Date of Acceptance: Feb 05, 2016

1. INTRODUCTION

Consider the hundreds, even thousands of products you can buy at your local grocery store. These items are often life-sustaining; almost every person on Earth nowadays relies on grocery stores for food and other household necessities. Because of this reliance, we have developed a certain expectation that our grocery stores will have what we need, when we need it, and for the most part, grocery stores around the world are successful in meeting this expectation. But, how do our grocery stores know which items to stock on their shelves, which items to direct our attention to in their marketing, which items to put together in packages or promotions? Increasingly, the answer to these types of questions is "data mining". More specifically, Association Rule modeling within the larger discipline of data mining has become a valuable tool for not only grocery stores, but many types of organizations, in determining patterns of behavior and relationships in large sets of data. This chapter will present a case example of how a supermarket can use simple modeling and data manipulation tools in Rapid Miner to create meaningful Association Rules.

Supermarkets are retail businesses, and modern retail businesses are predominantly man- aged through digital systems. Each transaction generates data, and as these data are aggregated, the resulting body of facts becomes extremely large. Just think about your own supermarket shopping behavior. When did you last go to the grocery store? How often do you go? How many products do you usually buy? How much do you spend each time you go? The answer to each of these questions leads to the generation of data-data which can be aggregated by product, product type or category, day of week, time of day - the list of possibilities is virtually endless. If we consider then not just your behavior, but the behavior of hundreds or thousands of customers who shop in a certain grocery store on a given day, we can see how quickly a dataset for a single store location can grow, not to mention a dataset for a chain or large conglomerate with multiple store locations. As items are passed across bar code scanners, money is collected, and receipts are handed to customers, the data continue to grow. These data can become a valuable asset to manage the store, but only if they can be analyzed quickly and accurately to inform store managers about customer preferences and choices. Since inventory databases, universal product bar codes and scanners, and other such supply chain management technologies have been around for years, the idea of using data to help manage retail operations is not new. However, more recently, the use of data mining to more thoroughly understand patterns of consumer behavior that affect retail operations has become more prevalent. In order to truly understand consumer behavior though, it is beneficial to understand both what they buy and who they are. Thus, in the past decade or so, we have seen an increase in the implementation of customer loyalty programs. You have probably seen these programs, and may even participate in them yourself. Generally, if you participate in such programs, you are given some form of reward, either a lower price on items in the store, or 'points' redeemable toward some future good or service. Airlines have been in the business of using such programs to encourage

customer loyalty for many years, with grocery and other

retail establishments adapting the concept to their operations more recently. But consider what you give when signing up for these programs. In order to receive the card which you subsequently use to gain the added benefit, you fill out a form. On this form, you give your name, gender, address, phone number, birth date, and perhaps any number of other personal characteristics. With this information, your grocer can go beyond traditional inventory management, and craft a much more personalized shopping experience for you. As we begin to examine how this might be accomplished, a comment about ethical behavior is in order. All organizations that collect, store, and analyze data have a responsibility to protect privacy, to guard against misuse and abuse, and to share data only within the constraints of fairly developed and disclosed policies. Aggregation of data, with the removal of personally identifiable information, is one way to ensure that peoples' privacy is protected. As a data miner, you bear a great responsibility to guard privacy and to behave ethically with the data you collect.

2. RELATED WORK

The association rule is one of the most widespread topics in the area of data mining[1].Since the concept of the association rules was put forward, many researchers have carried out a large number of studies on association rule mining problems, in which frequent item set mining algorithm is an important research direction. Among many algorithms, FP-Growth algorithm is the most famous, whose core is that all information of transactions are compressed to a FP-Tree in order to get the database mappings for frequent item set mining[2]. Afterwards, its conditional FP-Tree is constructed for each frequent item to mine frequent patterns.

Whereas, construction of a FP-Tree requires a lot of rams. When a transaction database is large to a certain extent, the running speed of the algorithm could be greatly reduced or the FP-Tree based on the ram cannot be constructed [3], which leads to the failure of mining.

2.1 FP GROWTH ALGORITHM

FP growth algorithm 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 ai in Q do {// Mining multipath FP-tree

9) Generate pattern β = ai \cup a with support = ai .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

Frequent itemsets are groups of items that often appear together in the data. It is important to know the basics of market-basket analysis for understanding frequent itemsets. The market-basket model of data is used to describe a common form of a many-to-many relationship between two kinds of objects. On the one hand, we have items, and on the other we have baskets, also called 'transactions'. The set of items is usually represented as set of attributes. Mostly these attributes are binominal. The transactions are usually each represented as examples of the Example Set. When an attribute value is 'true' in an example; it implies that the corresponding item is present in that transaction. Each transaction consists of a set of items (an itemset). Usually it is assumed that the number of items in a transaction is small, much smaller than the total number of items i.e. in most of the examples most of the attribute values are 'false'. The number of transactions is usually assumed to be very large i.e. the number of examples in the ExampleSet is assumed to be large. The frequent-itemsets problem is that of finding sets of items that appear together in at least a threshold ratio of transactions. This threshold is defined by the 'minimum support' criteria. The support of an itemset is the number of times that itemset appears in the ExampleSet divided by the total number of examples. The 'Purchases' data set at "Samples/data/Purchases" in the repository of RapidMiner is an example of how transactions data usually look like.

The discovery of frequent itemsets is often viewed as the discovery of 'association rules', although the latter is a more complex characterization of data, whose discovery depends fundamentally on the discovery of frequent itemsets. Association rules are derived from the frequent itemsets. The FP-Growth operator finds the frequent itemsets and operators like the Create Association Rules operator uses these frequent itemsets for calculating the association rules.

This operator calculates all frequent itemsets from an ExampleSet by building a FP-tree data structure on the transaction data base. This is a very compressed copy of the data which in many cases fits into main memory even for large data bases. All frequent itemsets are derived from this FP-tree. Many other frequent itemset mining algorithms also exist e.g. the Apriori algorithm. A major advantage of FP-Growth compared to Apriori is that it uses only 2 data scans and is therefore often applicable even on large data sets.

The given ExampleSet should contain only binominal attributes, i.e. nominal attributes with only two different values. If your ExampleSet does not satisfy this condition, you may use appropriate preprocessing operators to transform it into the required form. The discretization

operators can be used for changing the value of numerical attributes to nominal attributes. Then the Nominal to Binominal operator can be used for transforming nominal attributes into binominal attributes. The frequent itemsets are mined for the positive entries in your ExampleSet, i.e. for those nominal values which are defined as positive in your ExampleSet. If your data does not specify the positive entries correctly, you may set them using the positive value parameter. This only works if all your attributes contain this value.

This operator has two basic working modes:

 \cdot finding at least the specified number of itemsets with highest support without taking the 'min support' into account. This mode is available when the find min number of itemsets parameter is set to true. Then this operator finds the number of itemsets specified in the min number of itemsets parameter. The min support parameter is ignored in this case.

 \cdot finding all itemsets with a support larger than the specified minimum support. The minimum support is specified through the min support parameter. This mode is available when the find min number of itemsets parameter is set to false.

3. ASSOCIATION RULE MINING

One of the most common applications of ARM is market basket analysis (MBA) 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. 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 [4]. The aim of ARM analysis is to describe the most interesting patterns in an efficient manner [5].ARM analysis (also known as the market basket analysis (MBA)) is method of determining customer obtained patterns by mining association from retailer transactional database [6].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 [7].

3.1 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.



Figure 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	breads	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

Figure 2 Query results from an expanded dataset

ExampleSet (100151 examples, 0 special autobules, 12 legular autobules)									
Role	Name	Туре							
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	breads	binominal							
regular	produce	binominal							

Figure 3 the Meta data view

3.2 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.



Figure 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). When min support=0.75 and min support=0.5 Comparative Study 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.



Figure 5 Comparative Study of FP-Growth to our data mining process.

As we can see, the operator found frequencies for most items individually, and began to find frequencies between items as well. Although the screen capture does not show all 32 item sets that were found, it if did, you would be able to see that the final set found contains four products that appear to be associated with one another: juices, meats, frozen foods, and produce. There are a number of three-product combinations, and even more two-product sets. The Support attribute seen in Figure 6 indicates the number of observations in the dataset where the single or paired attributes was found; in other words, out of the 108,131.

When min support=0.5

Association Rules 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: [produce] --> [juices] (confidence: 0.781) 0.716)[meats, snack_foods] --> [juices] (confidence: 0.782) [juices, meats, frozen_foods] --> [produce] (confidence: [meats, snack_foods] --> [produce] (confidence: 0.783) 0.722) [meats, frozen_foods] --> [produce] (confidence: 0.725) [meats, beer_wine_spirits] --> [produce] (confidence: 0.784)[frozen_foods, produce] --> [meats] (confidence: 0.730) [produce, snack_foods] --> [juices] (confidence: 0.787) [juices, desserts] --> [frozen_foods] (confidence: 0.732) [desserts] --> [frozen_foods] (confidence: 0.737) [snack_foods] --> [juices] (confidence: 0.787) [juices, meats] --> [produce] (confidence: 0.739) [produce, paper_goods] --> [meats] (confidence: 0.788) [juices, produce] --> [meats] (confidence: 0.744) When min support=0.75 Association Rules as follows: [meats] --> [produce] (confidence: 0.744) [produce] --> [meats] (confidence: 0.752) [produce] --> [meats] (confidence: 0.752) [meats, frozen_foods, produce] --> [juices] (confidence: 0.760)[meats, frozen foods, produce] --> [juices] (confidence: 0.760)[paper goods] --> [juices] (confidence: 0.760) [paper_goods] --> [juices] (confidence: 0.760) [beer_wine_spirits] --> [juices] (confidence: 0.762) [beer_wine_spirits] --> [juices] (confidence: 0.762) [meats, frozen_foods] --> [juices] (confidence: 0.763) [meats, frozen_foods] --> [juices] (confidence: 0.763) [meats, beer_wine_spirits] --> [juices] (confidence: 0.764)[meats, beer_wine_spirits] --> [juices] (confidence: 0.764)[meats, paper_goods] --> [produce] (confidence: 0.766) [meats, paper_goods] --> [produce] (confidence: 0.766) [produce, beer_wine_spirits] --> [juices] (confidence: 0.767)[produce, beer_wine_spirits] --> [juices] (confidence: 0.767) [frozen foods] --> [juices] (confidence: 0.767) [frozen_foods, desserts] --> [juices] (confidence: 0.770) [frozen_foods] --> [juices] (confidence: 0.767) [frozen foods, desserts] --> [juices] (confidence: 0.770) [produce, snack foods] --> [meats] (confidence: 0.771) [produce, snack foods] --> [meats] (confidence: 0.771) [meats, produce] --> [juices] (confidence: 0.772) [meats, produce] --> [juices] (confidence: 0.772) [meats, paper_goods] --> [juices] (confidence: 0.772) [meats, paper_goods] --> [juices] (confidence: 0.772) [frozen_foods, produce] --> [juices] (confidence: 0.774) [frozen_foods, produce] --> [juices] (confidence: 0.774) [produce, beer_wine_spirits] --> [meats] (confidence: 0.775)[produce, beer wine spirits] --> [meats] (confidence: 0.775)[desserts] --> [juices] (confidence: 0.775) [desserts] --> [juices] (confidence: 0.775) [meats] --> [juices] (confidence: 0.777) [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)

4 CONCLUSIONS

We have found that juice products are relatively strongly connected to essentially every other product category in our grocery store, but what can we do with this information? Perhaps we already know, through daily experience, that we sell a lot of juice products, in which case this particular data mining model is of little help to us. But perhaps we might not have realized, without this model, just how pervasive juice products are in our product sales. As grocery store managers, we may begin to design product promotions which pair juice products with other strongly associated products in order to boost sales. We may go back and lower our confidence percentage a bit more, to see if other product categories emerge as the next most common conclusions (e.g., frozen foods and produce both have associations above 70% confidence). Or we may decide that we need more detail about specifically what juice products are frequently sold with other items, so we may choose to go back to the data extraction and preparation phase to group juice products into more specific attributes across our 108,131 receipts, to see if we might find even better clarity about what products our customers are most frequently buying together.

Authors



ACKNOWLEDGEMENTS:

A Project Supported by Scientific Research Fund of Hunan Provincial Education Department (15A043)

REFERENCES:

[1] J. S. Park, M. S. Chen, P. S. Yu. An effective Hash-based algorithm for mining association rules. In: Proc of 1995 ACM-SIGMOD Int'l Conf on Management of Data. San Jose. CA: ACM Press. 1995:175-186

[2] J. Han, J. Pei, Y. Yin. Mining frequent patterns without candidate generation. In: Proc of 2000 ACM-SIGMOD Int'l Conf on Management of Data. Dallas. TX: ACM Press. 2000:1-12

[3] X. P. Liu. Research and application of association rule mining algorithm based on FP-Growth algorithm. Hunan University, Hunan, 2006

[4] Chen, M., & Lin, C. A data mining approach to product assortment and shelf space allocation. Expert Systems with Applications, 2007,32,976-986.

[5] Tan, P., Steinbach. & Kumar, V. Introduction to data mining.Boston: Pearson Education, 2006

[6] Tang, K., Chen, Y., & Hu, H. (). Context-based market basket analysis in a multiple-store environment. DecisionSupportSystems, 2008, 45:150-163.

[7] Agrawal, R., Imielinski, T., & Swami, A.. Mining association rules between sets of items in large databases.ACM SIGMOD Conference, WashingtonDC, USA, 1993

< Zhi-hang Tang >, <1974-08-08>, <hunan, China> Current position, Doctor of Hunan Institute of Engineering University studies: control theory and control engineering in donghua University Scientific interest: intelligent decision and knowledge management Publications <number or main>: 30 Papers Experience: Zhihang TANG was born in Shaoyang, China, in 1974. He earned the M.S. degrees in control theory and control engineering from zhejiang University of techonlogy, in 2003 and Ph.D. from donghua University China in 2009. At the same time ,he is a teacher in department of computer and communication, Hunan Institute of Engineering(Xiangtan, China) from 2003.Chaired the 49th China Postdoctoral Science Foundation grant, presided over science and technology projects in Hunan Province in 2010, presided over the Education Department of Hunan Province in 2010 Outstanding Youth Project, as the first author more than 30 papers were published.His current research interests include intelligent decision and knowledge management.