

CHAPTER-14

Mining Association Rules in Large Databases

14.1 Introduction

14.2 Association Rule mining

14.3 Market Basket Analysis: A Motivating example for Association Rule Mining

14.4 Basic Concepts

14.5 Association Rule Mining: A Road map

14.6 Mining single-dimensional Association Rules from Transactional Databases

14.7 Generating Association Rules from Frequent Itemsets

14.8 Iceberg Queries

14.9 Review Questions

14.10 References

14. Mining Association Rules in Large Databases

14.1 Introduction

Association rule mining finds interesting association or correlation relationships among a large set of data items. With massive amounts of data continuously being collected and stored, many industries are becoming interested in mining association. Huge amounts of business transaction records can help in many business decision making processes, such as catalog design, cross-marketing, and loss-leader analysis.

A typical example of association rule mining is market basket analysis. This process analyzes customer buying habits by finding associations between the different items that customers place in their “shopping baskets”. The discovery of such associations can help retailers develop marketing strategies by gaining insight into which items are frequently purchased together by customers. For instance, if customers are buying milk, how likely are they to also buy bread (and what kind of bread) on the same trip to the supermarket? Such information can lead to increased sales by helping retailers do selective marketing and plan their shelf space. For example, placing milk and bread within close proximity may further encourage the sale of these items together within single visits to the store.

How can we find association rules from large amounts of data, where the data are either transactional or relational? Which association rules are the most interesting? How

can we help or guide the mining procedure to discover interesting associations? What language constructs are useful in defining a data mining query language for association rule mining.

14.2 Association Rule mining

Association rule mining searches for interesting relationships among items in a given data set. This section provides an introduction to association rule mining. We begin by presenting an example of market basket analysis, the earliest form of association rule mining. The basic concepts of mining associations are given and we present a road map to the different kinds of association rules that can be mined.

14.3 Market Basket Analysis: A Motivating example for Association Rule Mining

Suppose, as manager of an ABC company branch, you would like to learn more about the buying habits of your customers. Specifically, you wonder, "Which groups or sets of items are customers likely to purchase on a given trip to the store?" To answer your question, market basket analysis may be performed on the retail data of customer transactions at your store. The results may be used to plan marketing or advertising strategies, as well as catalog design. For instance, market basket analysis may help managers design different store layouts. In one strategy, items that are frequently purchased together can be placed in close proximity in order to further encourage the sale of such items together. If customers who purchase computers also tend to buy

Financial management software at the same time, then placing the hardware display close to the software display may help to increase the sales of both these items. In an alternative strategy, placing hardware and software at opposite ends of the store may entice customers who purchase such items to pick up other items along the way. For instance, after deciding on an expensive computer, a customer may observe security systems for sale while heading towards the software display to purchase financial management software and may decide to purchase a home security system as well. Market basket analysis can also help retailers to plan which items to put on sale at reduced prices. If customers tend to purchase computers and printers together, then having a sale on printers may encourage the sale of printers as well as computers.

If we think of the universe as the set of items available at the store, then each item has a Boolean variable representing the presence or absence of that item. A Boolean vector of values assigned to these variables can then represent each basket. The Boolean vectors can be analyzed for buying patterns that reflect items that are frequently associated or purchased together. These patterns can be represented in the form of association rules. For example, the information that customers who purchase computers also tend to buy financial management software at the same time is represented in the associated rule below:

Computer=> financial_management_software

[support = 2%,confidence == 60%]

Rule support and confidence are two measures of rule interesting that were described earlier. They respectively reflect the usefulness and certainty of discovered rules. A support of 2% for Association Rule means that 2% of all the transaction under analysis show that computer and financial management software are purchased a computer also bought the software. Typically, association rules are considered interesting if they satisfy both a minimum support thresholds and a minimum confidence threshold. Users or domain experts can set such thresholds.

14.4 Basic Concepts

Let $\tau = \{i_1, i_2, \dots, i_m\}$ be a set of items. Let D , the task-relevant data, be a set of database transactions where each transaction T is a set of items such that $T \subseteq \tau$ each transaction is associated with an identifier, called TID. Let A be a set of items. A transaction T is said to contain A if and only if $A \subseteq T$. An association rule is an implication of the form $A \Rightarrow B$, where $A \subseteq \tau$, $B \subseteq \tau$ and $A \cap B = \emptyset$. The rule $A \Rightarrow B$ holds in the transaction set D with support s , where s is the percentage of transaction in D that contains $A \cup B$ (i.e. both A and B). This is taken to be the probability, $P(A \cup B)$. The rule $A \Rightarrow B$ has confidence c in the transaction set D if c is the percentage of transactions in D containing A that also contain B . This is taken to be the conditional probability $P(B/A)$. That is

$$\text{Support}(A \Rightarrow B) = P(A \cup B)$$

$$\text{Confidence}(A \Rightarrow B) = P(B/A)$$

Rules that satisfy both a minimum support threshold (min_sup) and a minimum confidence threshold (min_conf) are called strong. By convention, we write support and confidence value so as to occur between 0% and 100% rather than 0 to 1.0.

A set of items is referred to as an itemset. An itemset that contains k items is a k -itemset. The set $\{\text{computer, financial_management_software}\}$ is a 2-itemset. The occurrence frequency of an itemset is the number of transactions that contain the itemset. This is also known, simply, as the frequency, support count or count of the itemset. An itemset satisfies minimum support if the occurrence frequency of the itemset is greater than or equal to the product of min_sup and the total number of transactions in D . The number of transactions required for the itemset to satisfy minimum support is

therefore referred to as the minimum support count. If an itemset satisfies minimum support, then it is a frequent itemset. The set of frequent K-itemsets is commonly denoted by L_K .

"How are association rules mined from large databases?" Association rule mining is a two-step process:

1. Find all frequent itemsets: By definition, each of these itemsets will occur at least as frequently as a pre-determined minimum support count.

2. Generate strong association rules from the frequent itemsets: By definition, these rules must satisfy minimum support and minimum confidence. Additional interestingness measures can be applied, if desired. The second step is the easiest, of the two. The overall performance of mining association rules is determined by the first step.

14.5 Association Rule Mining: A Road map

Market basket analysis is just one form of association rule mining, in fact, there are many kinds of association rules. Association rules can be classified in various ways, based on the following criteria:

Based on the types of values handled in the rule; if a rule concerns associations between the presence or absence of items, it is a Boolean association rule. For example, the rule above is a Boolean association rule obtained from market basket analysis.

If a rule describes associations between quantitative items or attributes, then it is a quantitative association rule. In these rules, quantitative values for items or attributes are partitioned into intervals. The following rule is an example of a quantitative association rule, where X is a variable representing a customer:

age(X,"30.....39")^income(X,"42K.....48")implies buys(X,high resolutionTV)

Note that the quantitative attributes,age and income,have been discretized,

Based on dimensions in the data: If the items or attributes in an association rule reference only one dimension,then it is a single-dimensional association rule,Note that above rule could be rewritten as buys(X,"computer") implies buys(X,"financial_management_software")The first rule above is a single-dimensional association rule since it refers to only one dimension,buys.If a rule references two or more dimensions,such as the dimensions buys,time_of_transaction,and customer_category,then it is a multidimensional association rule.

BASED on the levels of abstraction in the rule set:Some method for association rule mining can find rules at differing levels of abstraction. For example, suppose that a set of association rule mined includes the following rules:

age(x,"30...39") buys(x,"laptop computer")

age(x,"30...39") buys{x,"computer"}

In above rules the items bought are referenced at different levels of abstraction. (e.g.,"computer" is a higher-level abstraction of "laptop computer"). We refer to the rule set mined as consisting of multilevel association rules. If, instead,the rules within a given set do not reference items or attributes at different levels of abstraction, then then the set contains single-level association rules.

- **Based on various extensions to association mining:**Association mining can be extended to correlation analysis, where the absence or presence of correlated items can be identified. It can also be extended to mining maxpatterns (i.e.,maximal frequent patterns) and frequent closed itemsets. A **maxpattern** is a frequent pattern,p,such that any proper sub pattern of p is not frequent. A frequent closed itemset is a frequent closed itemset where an itemset c is closed if there exists no proper superset of c,c' such that every transaction containing c also contains c'. Maxpatterns and frequent closed itemset can be used to substantially reduce the number of frequent itemsets generated in mining.

14.6 Mining single-dimensional Association Rules from Transactional Databases:

In this section, you will learn methods for mining the simplest form of association rules—single-dimensional, single-level. Boolean association rules, such as those discussed for market basket analysis later. We begin by presenting a priori, a basic algorithm for improved efficiency and scalability are presented. Then methods for mining association rules that, unlike a priori, do not involve the generation of "candidate" frequent itemsets. The last section describes how principles from a priori can be applied to improve the efficiency of answering iceberg queries, which are common in market basket analysis.

The a priori Algorithm: Finding Frequent Itemsets Using Candidate Generation

A priori is an influential algorithm for mining frequent itemsets for boolean association rules. The name of the algorithm is based on the fact that the algorithm uses prior knowledge of frequent itemset properties, as we shall see below. A priori employs an iterative approach known as a level-wise search, where f -itemsets are used to explore $(k+1)$ -itemsets. First, the set of frequent 1-itemsets is found. This finding of each L_k requires one full scan of the database.

To improve the efficiency of the level-wise generation of frequent itemsets, an important property called the a priori property, presented below, is used to reduce the search space. We will first describe this property, and then show an example illustrating its use.

In order to use the a priori property, all nonempty subsets of a frequent itemset must also be frequent. This property is based on the following observation. By definition, if an itemset (i.e., $I \cup A$) cannot occur more frequently than I . Therefore, $I \cup A$ is not frequent either, that is, $P(I \cup A) < \min_sup$.

This property belongs to a special category of properties called anti-monotone in the sense that if a set cannot pass a test, all of its supersets will fail the same test as well. It is called anti-monotone because the property is monotonic in the context of failing a test.

To understand how the a priori property is used, let us look at how L_{k-1} is used to find L_k . A two-step process is followed, consisting of **join** and **prune** actions.

1. The join step: To find L_k , a set of candidate k -itemsets is generated by joining L_{k-1} with itself. This set of candidates is denoted C_k . Let I_1 and I_2 be itemsets in L_{k-1} ; The notation $I_1[j]$ refers to the j th item in I_1 ; (e.g., $I_1[k-2]$ refers to the second to the last item in I_1). By convention, a priori method assumes that items within a transaction or itemset are sorted in lexicographic order. The join, $L_{k-1} \times L_{k-1}$, is performed, where members of L_{k-1} are joinable if their first $(k-2)$ items are in common. That is, members I_1 and I_2 of L_{k-1} are joined if $(I_1[1]=I_2[1] \wedge I_1[k-2]=I_2[k-2] \wedge I_1[k-1] < I_2[k-2])$. The conditional $I_1[k-1] < I_2[k-2]$, simply ensures that no duplicates are generated. The resulting itemset formed by joining I_1 and I_2 is $I_1[I_1]I_2[k-1]$.

2. The prune step: C_k is a superset of L_k that is, its members may or may not be frequent, but all of the frequent k -itemsets are included in C_k . A scan of the database to determine the count of each candidate in C_k would result in the determination of L_k (i.e., all candidates having a count no less than the maximum support count are frequent by definition, and therefore belong to L_k). C_k , however, can be huge, and so this could involve heavy computation. To reduce the size of C_k , the a priori property is used as follows. Any $(k-1)$ -itemset that is not frequent cannot be a subset of a frequent k -itemset. Hence, if any $(k-1)$ -subset of a candidate k -itemset is not in L_{k-1} then the candidate cannot be frequent either and so can be removed from C_k . This subset testing can be done quickly by maintaining a hash tree of all frequent itemsets.

10.7 Generating Association Rules from Frequent Itemsets

Once the frequent itemsets from transactions in a database D have been found it is straightforward to generate strong association rules from them (where strong association rules satisfy both minimum support and minimum confidence). This can be done using the following equation for confidence, where the conditional probability is expressed in terms of itemset support count;

$$\text{confidence}(A \text{ implies } B) = P(B|A) = \frac{\text{support_count}\{A \cup B\}}{\text{support_count}\{A\}}$$

where $\text{support_count}(A \cup B)$ is the number of transactions containing the itemsets $A \cup B$, and $\text{support_count}(A)$ is the number of transactions containing the itemset A . Based on this equation, association rules can be generated as follows:

- For each frequent itemset l , generate all nonempty subsets of l .
- For every nonempty subset s of l , output the rule, $s \text{ implies } (l - s)$.

Since the rules are generated from frequent itemsets, each one automatically satisfies minimum support. Frequent itemsets can be stored ahead of time in hash tables along with their counts so that they can be accessed quickly.

10.8 Iceberg Queries

The Apriori algorithm can be used to improve the efficiency of answering ice-berg queries. Iceberg queries are commonly used in data mining, particularly for market basket analysis. An iceberg query computes an aggregate function over an attribute or set of attributes in order to find aggregate values above some specified threshold:

Given a relation R with attributes a_1, a_2, \dots, a_n and b , and an aggregate function, agg-f , an iceberg query is of the form

Given the large quantity of input data tuples, the number of tuples that will satisfy the threshold in the having clause is relatively small. The output result is seen as the "tip of the iceberg," where the "iceberg" is the set of input data.

14.9 Review Questions

- 1 Explain about Association Rule mining
- 2 Explain about Association Rule Mining: A Road map
- 3 Discuss Mining single-dimensional Association Rules from Transactional Databases
- 4 How can we Generate Association Rules from Frequent Itemsets
- 5 Explain about Iceberg Queries

14.10 References

- [1]. Data Mining Techniques, Arun K. Pujari 1st Edition
- [2]. Data Warehousing, Data Mining and OLAP, Alex Berson, Smith, J. Stephen
- [3]. Data Mining Concepts and Techniques, Jiawei Han and Micheline Kamber
- [4]. Data Mining Introductory and Advanced Topics, Margaret H. Dunham, PEA
- [5]. The Data Warehouse Lifecycle Toolkit, Ralph Kimball, Wiley Student Edition

