

Reasoning about Frequent Patterns with Negation

Marzena Kryszkiewicz

Warsaw University of Technology, Poland

INTRODUCTION

Discovering frequent patterns in large databases is an important data mining problem. The problem was introduced in (Agrawal, Imielinski, & Swami, 1993) for a sales transaction database. Frequent patterns were defined there as sets of items that are purchased together frequently. Frequent patterns are commonly used for building association rules. For example, an association rule may state that 80% of customers who buy fish also buy white wine. This rule is derivable from the fact that fish occurs in 5% of sales transactions and set {fish, white wine} occurs in 4% of transactions. Patterns and association rules can be generalized by admitting negation. A sample association rule with negation could state that 75% of customers who buy coke also buy chips and neither beer nor milk. The knowledge of this kind is important not only for sales managers, but also in medical areas (Tsumoto, 2002). Admitting negation in patterns usually results in an abundance of mined patterns, which makes analysis of the discovered knowledge infeasible. It is thus preferable to discover and store a possibly small fraction of patterns, from which one can derive all other significant patterns when required. In this chapter, we introduce first lossless representations of frequent patterns with negation.

BACKGROUND

Let us analyze sample transactional database \mathcal{D} presented in Table 1, which we will use throughout the chapter. Each row in this database reports items that were purchased by a customer during a single visit to a supermarket.

As follows from Table 1, items a and b were purchased together in four transactions. The number of transactions in which set of items $\{x_1, \dots, x_n\}$ occurs is called its *support* and denoted by $\text{sup}(\{x_1, \dots, x_n\})$. A set of items is called a *frequent pattern* if its support exceeds a user-specified threshold (minSup). Otherwise, it is called an *infrequent pattern*. In the remainder of the chapter, we assume $\text{minSup} = 1$. One can discover 27 frequent patterns from \mathcal{D} , which we list in Figure 1.

Table 1. Sample database \mathcal{D}

Id	Transaction
T_1	$\{abce\}$
T_2	$\{abcef\}$
T_3	$\{abch\}$
T_4	$\{abe\}$
T_5	$\{acfh\}$
T_6	$\{bef\}$
T_7	$\{h\}$
T_8	$\{af\}$

One can easily note that the support of a pattern never exceeds the supports of its subsets. Hence, subsets of a frequent pattern are also frequent, and supersets of an infrequent pattern are infrequent.

Aside from searching for only statistically significant sets of items, one may be interested in identifying frequent cases when purchase of some items (presence of some symptoms) excludes purchase of other items (presence of other symptoms). Pattern consisting of items x_1, \dots, x_m and negations of items x_{m+1}, \dots, x_n will be denoted by $\{x_1, \dots, x_m, \neg x_{m+1}, \dots, \neg x_n\}$. The *support of pattern* $\{x_1, \dots, x_m, \neg x_{m+1}, \dots, \neg x_n\}$ is defined as the number of transactions in which all items in set $\{x_1, \dots, x_m\}$ occur and no item in set $\{x_{m+1}, \dots, x_n\}$ occurs. In particular, $\{a(-b)\}$ is supported by two transactions in \mathcal{D} , while $\{a(-b)(-c)\}$ is supported by one transaction. Hence, $\{a(-b)\}$ is frequent, while $\{a(-b)(-c)\}$ is infrequent.

From now on, we will say that X is a *positive pattern*, if X does not contain any negated item. Otherwise, X is called a *pattern with negation*. A pattern obtained from pattern X by negating an arbitrary number of items in X is called a *variation of X* . For example, $\{ab\}$ has four distinct variations (including itself): $\{ab\}$, $\{a(-b)\}$, $\{(-a)b\}$, $\{(-a)(-b)\}$.

One can discover 109 frequent patterns in \mathcal{D} , 27 of which are positive, and 82 of which have negated items.

Figure 1. Frequent positive patterns discovered from database \mathcal{D} . Values provided in square brackets in the subscript denote supports of patterns.

$$\begin{array}{c}
 \{abce\}_{[2]} \\
 \{abc\}_{[3]} \quad \{abe\}_{[3]} \quad \{ace\}_{[2]} \quad \{acf\}_{[2]} \quad \{ach\}_{[2]} \quad \{bce\}_{[2]} \quad \{bef\}_{[2]} \\
 \{ab\}_{[4]} \quad \{ac\}_{[4]} \quad \{ae\}_{[3]} \quad \{af\}_{[3]} \quad \{ah\}_{[2]} \quad \{bc\}_{[3]} \quad \{be\}_{[4]} \quad \{bf\}_{[2]} \quad \{ce\}_{[2]} \quad \{cf\}_{[2]} \quad \{ch\}_{[2]} \quad \{ef\}_{[2]} \\
 \{a\}_{[6]} \quad \{b\}_{[5]} \quad \{c\}_{[4]} \quad \{e\}_{[4]} \quad \{f\}_{[4]} \quad \{h\}_{[3]} \\
 \emptyset_{[8]}
 \end{array}$$

In practice, the number of frequent patterns with negation is by orders of magnitude greater than the number of frequent positive patterns.

A first trial to solve the problem of large number of frequent patterns with negation was undertaken by Toivonen (1996), who proposed a method for using supports of positive patterns to derive supports of patterns with negation. The method is based on the observation that for any pattern X and any item x , the number of transactions in which X occurs is the sum of the number of transactions in which X occurs with x and the number of transactions in which X occurs without x . In other words, $\text{sup}(X) = \text{sup}(X \cup \{x\}) + \text{sup}(X \cup \{-x\})$, or $\text{sup}(X \cup \{-x\}) = \text{sup}(X) - \text{sup}(X \cup \{x\})$ (Mannila & Toivonen, 1996). Multiple usage of this property enables determination of the supports of patterns with an arbitrary number of negated items based on the supports of positive patterns. For example, the support of pattern $\{a(-b)(-c)\}$, which has two negated items, can be calculated as follows: $\text{sup}(\{a(-b)(-c)\}) = \text{sup}(\{a(-b)\}) - \text{sup}(\{a(-b)c\})$. Thus, the task of calculating the support of $\{a(-b)(-c)\}$, which has two negated items, becomes a task of calculating the supports of patterns $\{a(-b)\}$ and $\{a(-b)c\}$, each of which contains only one negated item. We note that $\text{sup}(\{a(-b)\}) = \text{sup}(\{a\}) - \text{sup}(\{ab\})$, and $\text{sup}(\{a(-b)c\}) = \text{sup}(\{ac\}) - \text{sup}(\{abc\})$. Eventually, we obtain: $\text{sup}(\{a(-b)(-c)\}) = \text{sup}(\{a\}) - \text{sup}(\{ab\}) - \text{sup}(\{ac\}) + \text{sup}(\{abc\})$. The support of $\{a(-b)(-c)\}$ is hence determinable from the supports of $\{abc\}$ and its proper subsets.

It was proved in Toivonen (1996) that for any pattern with negation its support is determinable from the supports of positive patterns. Nevertheless, the knowledge of the supports of only frequent patterns may be insufficient to derive the supports of all frequent patterns with negation (Boulicaut, Bykowski, & Jeudy, 2000), which we show beneath.

Let us try to calculate the support of pattern $\{bef(-h)\}$: $\text{sup}(\{bef(-h)\}) = \text{sup}(\{bef\}) - \text{sup}(\{befh\})$. Pattern $\{bef\}$ is frequent and its support equals 2 (see Figure 1). To the

contrary, $\{befh\}$ is not frequent, so its support does not exceed minSup , which equals 1. Hence, $1 \leq \text{sup}(\{bef(-h)\}) \leq 2$. The obtained result is not sufficient to determine if $\{bef(-h)\}$ is frequent.

The problem of large amount of mined frequent patterns is widely recognized. Within the last five years, a number of lossless representations of frequent positive patterns have been proposed. Frequent closed itemsets were introduced in (Pasquier et al., 1999); the generators representation was introduced in (Kryszkiewicz, 2001). Other lossless representations are based on disjunction-free sets (Bykowski & Rigotti, 2001), disjunction-free generators (Kryszkiewicz, 2001), generalized disjunction-free generators (Kryszkiewicz & Gajek, 2002), generalized disjunction-free sets (Kryszkiewicz, 2003), non-derivable itemsets (Calders & Goethals, 2002), and k -free sets (Calders & Goethals, 2003). All these models allow distinguishing between frequent and infrequent positive patterns and enable determination of supports for all frequent positive patterns. Although the research on concise representations of frequent positive patterns is advanced, no model was offered in the literature to represent all frequent patterns with negation.

MAIN THRUST

We offer a *generalized disjunction-free literal set model* (GDFLR) as a concise lossless representation of all frequent positive patterns and all frequent patterns with negation. Without the need to access the database, GDFLR enables distinguishing between all frequent and infrequent patterns, and enables calculation of the supports for all frequent patterns.

GDFLR uses the mechanism of deriving supports of positive patterns that was proposed in Kryszkiewicz & Gajek (2002). Hence, we first recall this mechanism. Then we examine how to use it to derive the supports of patterns with negation and propose a respective naive representation of frequent patterns. Next we examine relationships

4 more pages are available in the full version of this document, which may be purchased using the "Add to Cart" button on the publisher's webpage: www.igi-global.com/chapter/reasoning-frequent-patterns-negation/10731

Related Content

Approximate Range Queries by Histograms in OLAP

Francesco Buccafurri and Gianluca Lax (2005). *Encyclopedia of Data Warehousing and Mining* (pp. 49-53).
www.irma-international.org/chapter/approximate-range-queries-histograms-olap/10564

Distributed Data Management of Daily Car Pooling Problems

Roberto Wolfier Calvo, Fabio de Luigi, Palle Haastrup and Vittorio Maniezzo (2005). *Encyclopedia of Data Warehousing and Mining* (pp. 408-412).
www.irma-international.org/chapter/distributed-data-management-daily-car/10632

Ontology Query Languages for Ontology-Based Databases: A Survey

Stéphane Jean, Yamine Aït Ameur and Guy Pierra (2010). *Data Warehousing Design and Advanced Engineering Applications: Methods for Complex Construction* (pp. 227-247).
www.irma-international.org/chapter/ontology-query-languages-ontology-based/36617

Clustering Techniques

Sheng Ma and Tao Li (2005). *Encyclopedia of Data Warehousing and Mining* (pp. 176-179).
www.irma-international.org/chapter/clustering-techniques/10588

Mining E-Mail Data

Steffen Bickel and Tobias Scheffer (2008). *Data Warehousing and Mining: Concepts, Methodologies, Tools, and Applications* (pp. 1454-1460).
www.irma-international.org/chapter/mining-mail-data/7709