

HYPOTHESIS TESTING APPROACH IN MINING STATISTICALLY SIGNIFICANT FREQUENT ITEMSETS

PREDRAG STANIŠIĆ* AND SAVO TOMOVIĆ*

* Faculty of Science
University of Montenegro
Džordža Vašingtona bb, Podgorica, Montenegro
e-mail: pedjas@ac.me, web page: <http://www.pmf.ac.me>

Key words: frequent itemset mining, apriori algorithm

Summary. In this paper we suggest method in which hypothesis testing and support measure, which is used in Apriori algorithm, gives statistically significant frequent itemsets. Also, we propose the original procedure for frequent itemsets generation, which is more efficient than the appropriate procedure of the well known Apriori algorithm. The correctness of the procedure is based on special structure called Rymon tree. For its implementation, we suggest modified sort-merge-join algorithm.

1 INTRODUCTION

Association rules have received lots of attention in data mining due to their many applications in marketing, advertising, inventory control, and numerous other areas. The motivation for discovering association rules has come from the requirement to analyze large amounts of supermarket basket data. A record in such data typically consists of the transaction unique identifier and the items bought in that transaction. Items can be different products which one can buy in supermarkets or on-line shops, or car equipment, or telecommunication companies' services etc.

A typical supermarket may well have several thousand items on its shelves. Clearly, the number of subsets of the set of items is immense. Even though a purchase by a customer involves a small subset of this set of items, the number of such subsets is very large. For example, even if we assume that no customer has more than five items in his shopping cart, there are $\sum_{i=1}^5 \binom{1000}{i}$ possible contents of this cart, which corresponds to the subsets having no more than five items of a set that has 10,000 items, and this is indeed a large number.

The supermarket is interested in identifying associations between item sets; for example, it may be interesting to know how many of the customers who bought bread and cheese also bought butter. This knowledge is important because if it turns out that many of the customers who bought bread and cheese also bought butter, the supermarket will place butter physically close to bread and cheese in order to stimulate the sales of butter. Of course, such a piece of knowledge is especially interesting when there is a substantial number of customers who buy all three items and a large fraction of those individuals who buy bread and cheese also buy butter.

For example, the rule $bread \wedge cheese \Rightarrow butter$ [support=20%, confidence=85%] represents facts:

- 20% of all transactions under analysis contain bread, cheese and butter;
- 85% of the customers who purchased bread and cheese also purchased butter.

The result of association analysis is strong association rules, which are rules satisfying a minimal support and minimal confidence threshold. The minimal support and the minimal confidence are input parameters for association analysis.

The problem of association rules mining can be decomposed into two sub-problems [1]:

- Discovering frequent or large itemsets. Frequent itemsets have support greater than the minimal support;
- Generating rules. The aim of this step is to derive rules with high confidence (strong rules) from frequent itemsets. For each frequent itemset a one finds all not empty subsets of a ; for each $a \subseteq I \wedge a \neq \emptyset$ one generates the rule $a \Rightarrow I$, if $\frac{support(I)}{support(a)} \geq mins$.

We do not consider the second sub-problem in this paper, because the overall performances of mining association rules are determined by the first step. Efficient algorithms for solving the second sub-problem are exposed in [12].

The paper is organized as follows. Section 2 provides formalization of frequent itemsets mining problem. Section 3 describes Apriori algorithm. Section 4 presents our candidate generation procedure. Section 5 explains how to use hypothesis testing to validate frequent itemsets.

2 PRELIMINARIES

Suppose that I is a finite set; we refer to the elements of I as items.

Definition 2.1 A transaction dataset on I is a function $T: \{1, \dots, n\} \rightarrow P(I)$. The set $T(k)$ is the k^{th} transaction of T . The numbers $1, \dots, n$ are the transaction identifiers (TIDs).

Given a transaction data set T on the set I , we would like to determine those subsets of I that occur often enough as values of T .

Definition 2.2 Let $T: 1, \dots, n \rightarrow P(I)$ be a transaction data set on a set of items I . The support count of a subset K of the set of items I in T is the number $suppcount_T(K)$ given by:

$$suppcount_T(K) = |\{k \mid 1 \leq k \leq n \wedge K \subseteq T(k)\}|. \quad (1)$$

The support of an item set K (in the following text instead of an item set K we will use an itemset K) is the number:

$$supp_T(K) = suppcount_T(K)/n. \quad (2)$$

The following rather straightforward statement is fundamental for the study of frequent itemsets.

Theorem 2.1 Let $T: 1, \dots, n \rightarrow P(I)$ be a transaction data set on a set of items I . If K and K' are two itemsets, then $K' \subseteq K$ implies $supp_T(K') \leq supp_T(K)$.

Proof. The previous theorem states that $supp_T$ for an itemset has the anti-monotone property. It means that support for an itemset never exceeds the support for its subsets. For proof, it is sufficient to note that every transaction that contains K also contains K' . The statement from the theorem follows immediately.

Definition 2.3 An itemset K is μ -frequent relative to the transaction data set T if $supp_T(K) \geq \mu$. We denote by F_T^μ the collection of all μ -frequent itemsets relative to the transaction data set T and by $F_{T,r}^\mu$ the collection of μ -frequent itemsets that contain r items for $r \geq 1$ (in the following text we will use r -itemset instead of itemset that contains r items).

Note that $F_T^\mu = \bigcup_{r \geq 1} F_{T,r}^\mu$. If μ and T are clear from the context, then we will omit either or both adornments from this notation.

3 APRIORI ALGORITHM

We will briefly present the Apriori algorithm, which is the most famous algorithm in association analysis. The Apriori algorithm iteratively generates frequent itemsets starting from frequent 1-itemsets to frequent itemsets of maximal size. Each iteration of the algorithm consists of two phases: candidate generation and support counting. We will briefly explain both phases in the iteration k .

In the candidate generation phase potentially frequent k -itemsets or candidate k -itemsets are generated. The anti-monotone property of the itemset support is used in this phase and it provides elimination or pruning of some candidate itemsets without calculating its actual support (candidate containing at least one not frequent subset is pruned immediately, before support counting phase, see Theorem 2.1).

The support counting phase consists of calculating support for all previously generated candidates. In the support counting phase, it is essential to efficiently determine if the candidates are contained in particular transaction, in order to increment their support. Because of that, the candidates are organized in hash tree [2]. The candidates, which have enough support, are termed as frequent itemsets.

The Apriori algorithm terminates when none of the frequent itemsets can be generated.

We will use $C_{T,k}$ to denote candidate k -itemsets which are generated in iteration k of the Apriori algorithm. Pseudocode for the Apriori algorithm comes next.

Algorithm: Apriori

Input: A transaction dataset T ; Minimal support threshold μ

Output: F_T^μ

Method:

1. $C_{T,1} = \{\{i\} \mid i \in I\}$
2. $F_{T,1}^\mu = \{K \in C_{T,1} \mid supp_T(K) \geq \mu\}$
3. **for** ($k=2$; $F_{T,k-1}^\mu \neq \emptyset$; $k++$)

$$C_{T,k} = \text{apriori_gen}(F_{T,k-1}^\mu)$$

```


$$F_{T,k}^\mu = \{K \in C_{T,k} \mid \text{supp}_T \mu\}$$

end for
4.  $F_T^\mu = \bigcup_{r=1} F_{T,r}^\mu$ 

```

The **apriori_gen** is procedure for candidate generation in the Apriori algorithm [2]. In the iteration k , it takes as argument $F_{T,k-1}^\mu$ - the set of all μ -frequent $(k-1)$ -itemsets (which is generated in the previous iteration) and returns the set $C_{T,k}$. The set $C_{T,k}$ contains candidate k -itemsets (potentially frequent k -itemsets) and accordingly it contains the set $F_{T,k}^\mu$ - the set of all μ -frequent k -itemsets. The procedure works as follows. First, in the join step, we join $F_{T,k-1}^\mu$ with itself:

```

INSERT INTO  $C_{T,k}$ 
SELECT  $R_1.item_1, R_1.item_2, \dots, R_1.item_{k-1}, R_2.item_{k-1}$ 
FROM  $F_{T,k-1}^\mu$  AS  $R_1, F_{T,k-1}^\mu$  AS  $R_2$ 
WHERE  $R_1.item_1 = R_2.item_1 \wedge \dots \wedge R_1.item_{k-2} = R_2.item_{k-2} \wedge R_1.item_{k-1} < R_2.item_{k-1}$ 

```

In the previous query the set $F_{T,k-1}^\mu$ is considered as a relational table with $(k-1)$ attributes: $item_1, item_2, \dots, item_{k-1}$.

The previous query requires $k-2$ equality comparisons and is implemented with nested-loop join algorithm in the original Apriori algorithm from [2]. The nested-loop join algorithm is expensive, since it examines every pair of tuples in two relations:

```

for each  $l_1 \in F_{T,k-1}^\mu$ 
  for each  $l_2 \in F_{T,k-1}^\mu$ 
    if  $l_1[1] = l_2[1] \wedge \dots \wedge l_1[k-2] = l_2[k-2] \wedge l_1[k-1] < l_2[k-1]$ 
       $new\_cand = \{l_1[1], l_1[2], \dots, l_1[k-1], l_2[k-1]\}$ 

```

Consider the cost of nested-loop join algorithm. The number of pairs of itemsets (tuples) to be considered is $|F_{T,k-1}^\mu| * |F_{T,k-1}^\mu|$, where $|F_{T,k-1}^\mu|$ denotes the number of itemsets in $F_{T,k-1}^\mu$. For each itemset from $F_{T,k-1}^\mu$, we have to perform a complete scan on $F_{T,k-1}^\mu$. In the worst case the buffer can hold only one block of $F_{T,k-1}^\mu$, and a total of $|F_{T,k-1}^\mu| * n_b + n_b$ block transfers would be required, where n_b denotes the number of blocks containing itemsets of $F_{T,k-1}^\mu$.

Next in the *prune* step, we delete all candidate itemsets $c \in C_{T,k}$ such that some $(k-1)$ -subset of c is not in $F_{T,k-1}^\mu$ (see Theorem 2.1.):

```

for each  $c \in C_{T,k}$ 
  for each  $(k-1)$ -subset  $s$  in  $c$ 
    if  $s \notin F_{T,k-1}^\mu$  delete  $c$  from  $C_{T,k}$ 

```

In the next section we will suggest more efficient procedure for candidate generation phase and an efficient algorithm for its implementation.

4 NEW CANDIDATE GENERATION PROCEDURE

We assume that any itemset K is kept sorted according to some relation $<$, where for all $x, y \in K$, $x < y$ means that the object x is in front of the object y . Also, we assume that all transactions in database T and all subsets of K are kept sorted in lexicographic order according to the relation $<$.

For candidate generation we suggest the original method by which the set $C_{T,k}$ is calculated by joining $F_{T,k-1}^\mu$ with $F_{T,k-2}^\mu$, in the iteration k , for $k \geq 3$. Candidate k -itemset is created from one frequent $(k-1)$ -itemset and one frequent $(k-2)$ -itemset in the following way. Let $X = \{x_1, \dots, x_{k-1}\} \in F_{T,k-1}^\mu$ and $Y = \{y_1, \dots, y_{k-2}\} \in F_{T,k-2}^\mu$. Itemsets X and Y are joined if and only if the following condition is satisfied:

$$x_i = y_i, (1 \leq i \leq k-3) \wedge x_{k-1} < y_{k-2} \quad (3)$$

producing the candidate k -itemset $\{x_1, \dots, x_{k-2}, x_{k-1}, y_{k-2}\}$.

We will prove the correctness of the suggested method. In the following text we will denote this method by $C_{T,k} = F_{T,k-1}^\mu \times F_{T,k-2}^\mu$. Let $I = i_1, \dots, i_n$ be a set of items that contains n elements. Denote by $G_I = (P(I), E)$ the Rymon tree (see Appendix) of $P(I)$. The root of the tree is \emptyset . A vertex $K = \{i_{p_1}, \dots, i_{p_k}\}$ with $i_{p_1} < i_{p_2} < \dots < i_{p_k}$ has $n - i_{p_k}$ children $K \cup j$, where $i_{p_k} < j \leq n$ (see Appendix). Let S_r be the collection of itemsets that have r elements. The next theorem suggest a technique for generating S_r starting from S_{r-1} and S_{r-2} .

Theorem 4.1 Let G be the Ryman tree of $P(I)$, where $I = i_1, \dots, i_n$. If $W \in S_r$, where $r \geq 3$, then there exists a unique pair of distinct sets $U \in S_{r-1}$ and $V \in S_{r-2}$ that has a common immediate ancestor $T \in S_{r-3}$ in G such that $U \cap V \in S_{r-3}$ and $W = U \cup V$.

Proof. Let u and v and p be the three elements of W that have the largest, the second-largest and the third-largest subscripts, respectively. Consider the sets $U = W - \{u\}$ and $V = W - \{v, p\}$. Note that $U \in S_{r-1}$ and $V \in S_{r-2}$. Moreover, $Z = U \cup V$ belongs to S_{r-3} because it consists of the first $r-3$ elements of W . Note that both U and V are descendants of Z and that $U \cup V = W$ (for $r=3$ we have $Z = \emptyset$).

The pair (U, V) is unique. Indeed, suppose that W can be obtained in the same manner from another pair of distinct sets $U_1 \in S_{r-1}$ and $V_1 \in S_{r-2}$ such that U_1 and V_1 are immediate descendants of a set $Z_1 \in S_{r-3}$. The definition of the Rymon tree G_I implies that $U_1 = Z_1 \cup \{i_m, i_q\}$ and $V_1 = Z_1 \cup \{i_y\}$, where the letters in Z_1 are indexed by a number smaller than $\min\{m, q, y\}$. Then, Z_1 consists of the first $r-3$ symbols of W , so $Z_1 = Z$. If $m < q < y$,

then m is the third-highest index of a symbol in W , q is the second-highest index of a symbol in W and y is the highest index of a symbol in W , so $U_1 = U$ and $V_1 = V$.

The following theorem directly proves correctness of our method $C_{T,k} = F_{T,k-1}^\mu \times F_{T,k-2}^\mu$.

Theorem 4.2 *Let T be a transaction data set on a set of items I and let $k \in N$ such that $k > 2$. If W is a μ -frequent itemset and $|W| = k$, then there exists a μ -frequent itemset Z and two itemsets $\{i_m, i_q\}$ and $\{i_y\}$ such that $|Z| = k - 3$, $Z \subseteq W$, $W = Z \cup \{i_m, i_q, i_y\}$ and both $Z \cup \{i_m, i_q\}$ and $Z \cup \{i_y\}$ are μ -frequent itemsets.*

Proof. If W is an itemset such that $|W| = k$, then we already know that W is the union of two subsets U and V of I such that $|U| = k - 1$, $|V| = k - 2$ and that $Z = U \cap V$ has $k - 3$ elements (it follows from Theorem 2). Since W is a μ -frequent itemset and Z , U , V are subsets of W , it follows that each of these sets is also a μ -frequent itemset (it follows from Theorem 2.1).

We have seen in the previous section than in the original Apriori algorithm from [2], the join procedure is suggested. It generates candidate k -itemset by joining two large $(k-1)$ -itemsets, if and only if they have first $(k-2)$ items in common. Because of that, each join operation requires $(k-2)$ equality comparisons. If a candidate k -itemset is generated by the method $C_{T,k} = F_{T,k-1}^\mu \times F_{T,k-2}^\mu$ for $k \geq 3$, it is enough $(k-3)$ equality comparisons to process.

The method $C_{T,k} = F_{T,k-1}^\mu \times F_{T,k-2}^\mu$ can be represented by the following SQL query:

```

INSERT INTO CT,k
SELECT R1.item1, ..., R1.itemk-1, R2.itemk-2
FROM FT,k-1μ AS R1, FT,k-2μ AS R2
WHERE R1.item1 = R2.item1 ∧ ... ∧ R1.itemk-3 = R2.itemk-3 ∧ R1.itemk-1 < R2.itemk-2

```

For the implementation of the join $C_{T,k} = F_{T,k-1}^\mu \times F_{T,k-2}^\mu$ we suggest a modification of sort-merge-join algorithm (note that $F_{T,k-1}^\mu$ and $F_{T,k-2}^\mu$ are sorted because of the way they are constructed and lexicographic order of itemsets).

By the original sort-merge-join algorithm [9], it is possible to compute natural joins and equi-joins. Let $r(R)$ and $s(S)$ be the relations and $R \cap S$ denote their common attributes. The algorithm keeps one pointer on the current position in relation $r(R)$ and another one pointer on the current position in relation $s(S)$. As the algorithm proceeds, the pointers move through the relations. It is supposed that the relations are sorted according to joining attributes, so tuples with the same values on the joining attributes are in consecutive order. Thereby, each tuple needs to be read only once, and, as a result, each relation is also read only once.

The number of blocks transfers is equal to the sum of the number of blocks in both sets $F_{T,k-1}^\mu$ and $F_{T,k-2}^\mu$, $n_{b1} + n_{b2}$. We have seen that nested-loop join requires $|F_{T,k-1}^\mu| * n_{b1} + n_{b1}$ block transfers and we can conclude that merge-join is more efficient ($|F_{T,k-1}^\mu| \gg n_{b2}$).

The modification of sort-merge-join algorithm we suggest refers to the elimination of restrictions that join must be natural or equi-join. First, we separate the condition (3):

$$x_i = y_i, 1 \leq i \leq k-3 \quad (4)$$

$$x_{k-1} < y_{k-2}. \quad (5)$$

Joining $C_{T,k} = F_{T,k-1}^\mu \times F_{T,k-2}^\mu$ is calculated according to the condition (4), in other words we compute natural join. For this, the described sort-merge-join algorithm is used, and our modification is: before $X = \{x_1, \dots, x_{k-1}\}$ and $Y = \{y_1, \dots, y_{k-2}\}$, for which $X \in F_{T,k-1}^\mu$ and $Y \in F_{T,k-2}^\mu$ and $x_i = y_i, 1 \leq i \leq k-3$ is true, are joined, we check if condition (5) is satisfied, and after that we generate candidate k -itemset $\{x_1, \dots, x_{k-2}, x_{k-1}, y_{k-2}\}$.

The pseudocode of **apriori_gen** function comes next.

FUNCTION **apriori_gen**($F_{T,k-1}^\mu, F_{T,k-2}^\mu$)

1. $i = 0$

2. $j = 0$

3. while $i \leq |F_{T,k-1}^\mu| \wedge j \leq |F_{T,k-2}^\mu|$

$iset_1 = F_{T,k-1}^\mu[i++]$

$S = \{iset_1\}$

done = false

while $done = false \wedge i \leq |F_{T,k-1}^\mu|$

$iset_{1a} = F_{T,k-1}^\mu[i++]$

if $iset_{1a}[w] = iset_1[w], 1 \leq w \leq k-2$ then

$S = S \cup \{iset_{1a}\}$

$i++$

else

done = true

end if

end while

$iset_2 = F_{T,k-2}^\mu[j]$

while $j \leq |F_{T,k-2}^\mu| \wedge iset_2[1, \dots, k-2] \prec iset_1[1, \dots, k-2]$

$iset_2 = F_{T,k-2}^\mu[j++]$

end while

while $j \leq |F_{T,k-2}^\mu| \wedge iset_1[w] = iset_2[w], 1 \leq w \leq k-2$

for each $s \in S$

if $iset_1[k-1] \prec iset_2[k-2]$ then

$c = \{iset_1[1], \dots, iset_1[k-1], iset_2[k-2]\}$

if c contains-not-frequent-subset then DELETE c

else

```

         $C_{T,k} = C_{T,k} \cup \{c\}$ 
    end if
end for
j++
 $iset_2 = F_{T,k-2}^\mu[j]$ 
end while
end while

```

We implemented the original Apriori [1] and the algorithm proposed here. Algorithms are implemented in C in order to evaluate its performances. Experiments are performed on PC with a CPU Intel(R) Core(TM)2 clock rate of 2.66GHz and with 2GB of RAM. Also, run time used here means the total execution time, i.e., the period between input and output instead of CPU time measured in the experiments in some literature. In experiments dataset which can be found on www.cs.uregina.ca is used. It contains 10000 binary transactions. The average length of transactions is 8.

Table 1 shows that the original Apriori algorithm from [1] is outperformed by the algorithm presented here.

<i>minsup</i> (%)	<i>Apriori</i>	<i>Rymon tree based implementation</i>
10.0	15.00	3.10
5.0	15.60	6.20
2.5	16.60	6.20
1.0	16.60	7.80
0.5	18.10	7.80

Table 1: Execution time in seconds

5 VALIDATING FREQUENT ITEMSETS

In this section we use hypothesis testing to validate frequent itemsets, which have been generated in the Apriori algorithm.

Hypothesis testing is a statistical inference procedure to determine whether a hypothesis should be accepted or rejected based on the evidence gathered from data. Examples of hypothesis tests include verifying the quality of patterns extracted by many data mining algorithms and validating the significance of the performance difference between two classification models.

In hypothesis testing, we are usually presented with two contrasting hypothesis, which are known, respectively, as the null hypothesis and the alternative hypothesis. The general procedure for hypothesis testing consists of the following four steps:

- Formulate the null and alternative hypotheses to be tested.
- Define a test statistic θ that determines whether the null hypothesis should be accepted or rejected. The probability distribution associated with the test statistic should be known.

- Compute the value of θ from the observed data. Use the knowledge of the probability distribution to determine a quantity known as p -value.

- Define a significance level α which controls the range of θ values in which the null hypothesis should be rejected. The range of values for θ is known as the rejection region.

Now, we will back to mining frequent itemsets and we will evaluate the quality of the discovered frequent itemsets from a statistical perspective. The criterion for judging whether the itemset is frequent depends on the support of the itemset. Recall that support measures the number of transactions in which the itemset is actually observed. An itemset X is considered frequent in the data set T , if $supp_T(X) > minsup$, where $minsup$ is a user-specified threshold.

The problem can be formulated into the hypothesis testing framework in the following way. To validate if the itemset X is frequent in the data set T , we need to decide whether to accept the null hypothesis, $H_0 : supp_T(X) = minsup$, or the alternative hypothesis $H_1 : supp_T(X) > minsup$. If the null hypothesis is rejected, then X is considered as frequent itemset. To perform the test, the probability distribution for $supp_T(X)$ must also be known.

Theorem 5.1 The measure $supp_T(X)$ for the itemset X in transaction data set T has the binomial distribution with mean $supp_T(X)$ and variance

$$\frac{supp_T(X) * (1 - supp_T(X))}{n}, \text{ where } n \text{ is the number of transactions in } T.$$

Proof. We will use measure $suppcount_T(X)$ and calculate mean and variance for it and later derive mean and variance for the measure $supp_T(X)$. The measure $suppcount_T(X) = X_n$ presents the number of transactions in T that contain itemset X , and $supp_T(X) = suppcount_T(X)/n$ (Definition 2.2).

The measure X_n is analogous to determining the number of heads that shows up when tossing n coins. Let us calculate $E(X_n)$ and $D(X_n)$.

Mean is $E(X_n) = n * p$, where p is the probability of success, which means (in our case) the itemset X appears in one transaction. According to Bernoulli law is:

$$\forall \varepsilon > 0, \lim_{N \rightarrow \infty} P\left\{ \left| \frac{X_n}{n} - p \right| \leq \varepsilon \right\} = 1. \quad (6)$$

Freely speaking, for large n (we work with large databases so n can be considered large), we can use relative frequency instead of probability. So, we now have:

$$E(X_n) = np \approx n \frac{X_n}{n} = X_n \quad (7)$$

For variance we compute:

$$D(X_n) = np(1-p) \approx n \frac{X_n}{n} \left(1 - \frac{X_n}{n}\right) = \frac{X_n(n - X_n)}{n} \quad (8)$$

Now we will compute $E(supp_T(X))$ and $D(supp_T(X))$. Recall that

$supp_T(X) = \frac{X_n}{n}$. We have:

$$E(supp_T(X)) = E\left(\frac{X_n}{n}\right) = \frac{1}{n}E(X_n) = \frac{X_n}{n} = supp_T(X). \quad (9)$$

$$D(supp_T(X)) = D\left(\frac{X_n}{n}\right) = \frac{1}{n^2}D(X_n) = \quad (10)$$

$$= \frac{1}{n^2} \frac{X_n(n - X_n)}{n} = \frac{1}{n} supp_T(X)(1 - supp_T(X)) \quad (11)$$

◇

The binomial distribution can be further approximated using a normal distribution if n is sufficiently large, which is typically the case in association analysis.

Regarding previous paragraph and Theorem 5.1, under the null hypothesis $supp_T(X)$ is assumed to be normally distributed with mean $minsup$ and variance $\frac{minsup(1 - minsup)}{n}$. To test whether the null hypothesis should be accepted or rejected, the following statistic can be used:

$$W_N = \frac{supp_T(X) - minsup}{\sqrt{\frac{minsup(1 - minsup)}{n}}}. \quad (12)$$

The previous statistic, according to the Central Limit Theorem, has the distribution $N(0,1)$. The statistic essentially measures the difference between the observed support $supp_T(X)$ and the $minsup$ threshold in units of standard deviation.

Let $N=10000$, $supp_T(X) = 0.11$, $minsup=0.1$ and $\alpha = 0.001$. The last parameter is the desired significance level. It controls Type 1 error which is rejecting the null hypothesis even though the hypothesis is true.

In the Apriori algorithm we compare $supp_T(X) = 0.11 > 0.1 = minsup$ and we declare X as frequent itemset. Is this validation procedure statistically correct?

Under the hypothesis H_1 statistics W_{10000} is positive and for rejection region we choose

$$R = \{(x_1, \dots, x_{10000}) \mid w_{10000} > k\}, k > 0.$$

Let us find k .

$$0.001 = P_{H_0} \{W_{10000} > k\}$$

$$P_{H_0} \{W_{10000} > k\} = 0.499$$

$$k = 3.09$$

Now we compute $w_{10000} = \frac{0.11-0.1}{\sqrt{\frac{0.1*(1-0.1)}{10000}}} = 3.33\dots$

We can see that $w_{10000} > k$, so we are in rejection region and H_1 is accepted, which means the itemset X is considered statistically significant.

6 APPENDIX

Rymon tree was introduced in [8] in order to provide a unified search-based framework for several problems in artificial intelligence; the Rymon tree is also useful for data mining algorithms.

In Definitions 6.1 and 6.2 we define necessary concepts and in Definition 6.3 we define the Rymon tree.

Definition 6.1. Let S be a set and let $d : S \rightarrow N$ be an injective function. The number $d(x)$ is the index of $x \in S$. If $P \subseteq S$, the view of P is the subset $view(d, P) = \{s \in S \mid d(s) > \max_{p \in P} d(p)\}$.

Definition 6.2. A collection of sets C is hereditary if $U \in C$ and $W \subseteq U$ implies $W \in C$.

Definition 6.3. Let C be a hereditary collection of subsets of a set S . The graph $G = (C, E)$ is a Rymon tree for C and the indexing function d if:

- the root of the G is \emptyset
- the children of a node P are the sets of the form $P \cup \{s\}$, where $s \in view(d, P)$

If $S = \{s_1, \dots, s_n\}$ and $d(s_i) = i$ for $1 \leq i \leq n$, we will omit the indexing function from the definition of the Rymon tree for $P(S)$.

A key property of a Rymon tree is stated next.

Theorem 6.1. Let G be a Rymon tree for a hereditary collection C of subsets of a set S and an indexing function d . Every set P of C occurs exactly once in the tree.

Proof. The argument is by induction on $p = |P|$. If $p = 0$, then P is the root of the tree and the theorem obviously holds.

Suppose that the theorem holds for sets having fewer than p elements, and let $P \in C$ be such that $|P| = p$. Since C is hereditary, every set of the form $P - \{x\}$ with $x \in P$ belongs to C and, by the inductive hypothesis, occurs exactly once in the tree.

Let z be the element of P that has the largest value of the index function d . Then $view(P - \{z\})$ contains z and P is a child of the vertex $P - \{z\}$. Since the parent of P is

unique, it follows that P occurs exactly once in the tree. \diamond

Note that in the Rymon tree of a collection of the form $P(S)$, the collection of sets of S_r that consists of sets located at distance r from the root denotes all subsets of the size r of S .

REFERENCES

- [1] R. Agrawal, R. Srikant, *Fast Algorithms for Mining Association Rules*, Proceedings of VLDB-94 (1994), 487-499.
- [2] F.P. Coenen, P. Leng, S. Ahmed, *T-Trees, Vertical Partitioning and Distributed Association Rule Mining*, Proceedings ICDM-2003 (2003), 513-516.
- [3] F.P. Coenen, P. Leng, S. Ahmed, *Data Structures for Association Rule Mining: T-trees and P-trees*, IEEE Transactions on Data and Knowledge Engineering, Vol. 16, No 6, (2004), 774-778.
- [4] F.P. Coenen, P. Leng, G. Goulbourne, *Tree Structures for Mining Association Rules*, Journal of Data Mining and Knowledge Discovery, Vol. 8, No. 1, (2004), 25-51.
- [5] G. Goulbourne, F.P. Coenen, P. Leng, *Algorithms for Computing Association Rules Using a Partial-Support Tree*, Journal of Knowledge-Based Systems, Vol. 13, (1999), 141-149.
- [6] G. Grahne, J. Zhu, *Efficiently Using Prefix-trees in Mining Frequent Itemsets*, Proceedings of the IEEE ICDM Workshop on Frequent Itemset Mining Implementations (2003).
- [7] J. Han, J. Pei, P.S. Yu, *Mining Frequent Patterns without Candidate Generation*, Proceedings of the ACM SIGMOD Conference on Management of Data, (2000), 1-12.
- [8] R. Rymon, *Search Through Systematic Set Enumeration*, Proceedings of 3rd International Conference on Principles of Knowledge Representation and Reasoning, (1992), 539-550.
- [9] Silberschatz, A., Korth, H. F., Sudarshan, S., *Database System Concepts*, Mc Graw Hill, New York (2006)
- [10] A. D. Simovici, C. Djeraba, *Mathematical Tools For Data Mining (Set Theory, Partial Orders, Combinatorics)*, Springer-Verlag, London, 2008.
- [11] P. Stanisic, S. Tomovic, *Apriori Multiple Algorithm for Mining Association Rules*, Information Technology and Control, Vol. 37, No. 4, (2008), 311-320.
- [12] P.N. Tan., M. Steinbach, V. Kumar, *Introduction to Data Mining*, Addison Wesley, 2006.