

Mining Weighted Association Rules for Fuzzy Quantitative Items

Attila Gyenesei

Turku Centre for Computer Science (TUCS)
University of Turku, Department of Computer Science,
Lemminkaisenkatu 14, FIN-20520 Turku, Finland
gyenesei@cs.utu.fi

Abstract. During the last ten years, data mining, also known as knowledge discovery in databases, has established its position as a prominent and important research area. Mining association rules is one of the important research problems in data mining. Many algorithms have been proposed to find association rules in large databases containing both categorical and quantitative attributes. We generalize this to the case where part of attributes are given weights to reflect their importance to the user.

In this paper, we introduce the problem of mining weighted quantitative association rules based on fuzzy approach. Using the fuzzy set concept, the discovered rules are more understandable to a human.

We propose two different definitions of weighted support: with and without normalization. In the normalized case, a subset of a frequent itemset may not be frequent, and we cannot generate candidate k -itemsets simply from the frequent $(k-1)$ -itemsets. We tackle this problem by using the concept of *z-potential frequent subset* for each candidate itemset.

We give an algorithm for mining such quantitative association rules. Finally, we describe the results of using this approach on a real-life dataset.

1 Introduction

The goal of data mining is to extract higher level information from an abundance of raw data. Mining association rules is one of the important research problems in data mining [1]. The problem of mining boolean association rules was first introduced in [2], and later broadened in [3], for the case of databases consisting of categorical attributes alone. Categorical association rules are rules where the events X and Y , on both sides of the rule, are instances of given categorical items. In this case, we wish to find all rules with confidence and support above user-defined thresholds (minconf and minsup). Several efficient algorithms for mining categorical association rules have been published (see [3], [4], [5] for just a few examples).

A variation of categorical association rules was recently introduced in [6]. Their new definition is based on the notion of weighted items to represent the importance of individual items.

The problem of mining quantitative association rules was introduced and an algorithm proposed in [7]. The algorithm finds the association rules by partitioning the attribute domain, combining adjacent partitions, and then transforming the problem into binary one. An example of a rule according to this definition would be: “10% of married people between age 50 and 70 have at least 2 cars”.

In [8], we showed a method to handle quantitative attributes using a fuzzy approach. We assigned each quantitative attribute several fuzzy sets which characterize it. Fuzzy sets provide a smooth transition between a member and non-member of a set. The fuzzy association rule is also easily understandable to a human because of the linguistic terms associated with the fuzzy sets. Using the fuzzy set concept, the above example could be rephrased e.g. “10% of married old people have several cars”.

In this paper, we introduce a new definition of the notion of weighted itemsets based on fuzzy set theory. In a marketing business, a manager may want to mine the association rules with more emphasis on some itemsets in mind, and less emphasis on other itemsets. For example, some itemsets may be more interesting for the company than others. This results in a generalized version of the quantitative association rule mining problem, which we call weighted quantitative association rule mining.

The paper is organized as follows. In the next section, we will present the definition of mining quantitative association rules using a fuzzy approach. Then we will introduce the problem of weighted quantitative association rules in Section 3. In Section 4, we give a new algorithm for this problem. In Section 5 the experimental results are reported, followed by a brief conclusion in Section 6.

2 Fuzzy Association Rules

In [8], an algorithm for mining quantitative association rules using a fuzzy approach was proposed. We summarize its definitions in what follows.

Let $I = \{i_1, i_2, \dots, i_m\}$ be the complete item set where each i_j ($1 \leq j \leq m$) denotes a categorical or quantitative (fuzzy) attribute. Suppose $f(i_j)$ represents the maximum number of categories (if i_j is categorical) or the maximum number of fuzzy sets (if i_j is fuzzy), and $d_{i_j}(l, v)$ represents the membership degree of v in the l^{th} category or fuzzy set of i_j . If i_j is categorical, $d_{i_j}(l, v) = 0$ or $d_{i_j}(l, v) = 1$. If i_j is fuzzy, $0 \leq d_{i_j}(l, v) \leq 1$.

Let $t = \{t.i_1, t.i_2, \dots, t.i_m\}$ be a transaction, where $t.i_j$, ($1 \leq j \leq m$) represents a value of the j^{th} attribute and can be mapped to

$$\{(l, d_{i_j}(l, t.i_j)) \mid \text{for all } l, 1 \leq l \leq f(i_j)\}.$$

Given a database $D = \{t_1, t_2, \dots, t_n\}$ with attributes I and the fuzzy sets associated with attributes in I , we want to find out some interesting, potentially useful regularities.

Definition 1. A *fuzzy association rule* is of the form

If $X = \{x_1, x_2, \dots, x_p\}$ is $A = \{a_1, a_2, \dots, a_p\}$ then $Y = \{y_1, y_2, \dots, y_q\}$ is $B = \{b_1, b_2, \dots, b_q\}$, where X, Y are itemsets and $a_i \in \{\text{fuzzy sets related to attribute } x_i\}$, $b_j \in \{\text{fuzzy sets related to attribute } y_j\}$

X and Y are ordered subsets of I and they are disjoint i.e. they share no common attributes. A and B contain the fuzzy sets associated with the corresponding attributes in X and Y . As in the binary association rule, “ X is A ” is called the antecedent of the rule while “ Y is B ” is called the consequent of the rule.

If a rule is interesting, it should have enough support and a high confidence value. We define the terms support and confidence as in [8]: the fuzzy support value is calculated by first summing all votes of each record with respect to the specified itemset, then dividing it by the total number of records. Each record contributes a vote which falls in $[0, 1]$. Therefore, a fuzzy support value reflects not only the number of records supporting the itemset, but also their degree of support.

Definition 2. The *fuzzy support value* of itemset $\langle X, A \rangle$ in transaction set D is

$$FS_{\langle X, A \rangle} = \frac{\sum_{t_i \in D} \prod_{x_j \in X} d_{x_j}(a_j, t_i \cdot x_j)}{|D|}$$

Definition 3. An itemset $\langle X, A \rangle$ is called a *frequent itemset* if its fuzzy support value is greater than or equal to the minimum support threshold.

We use the discovered frequent itemsets to generate all possible rules. If the union of antecedent $\langle X, A \rangle$ and consequent $\langle Y, B \rangle$ has enough support and the rule has high confidence, this rule will be considered as interesting.

When we obtain a frequent itemset $\langle Z, C \rangle$, we want to generate fuzzy association rules of the form, “If X is A then Y is B ”, where $X \subset Z$, $Y = Z - X$, $A \subset C$ and $B = C - A$. Having the frequent itemset, we know its support as well as the fact that all of its subsets will be also frequent.

Definition 4. The *fuzzy confidence value* of a rule is as follows:

$$FC_{\langle \langle X, A \rangle, \langle Y, B \rangle \rangle} = \frac{FS_{\langle Z, C \rangle}}{FS_{\langle X, A \rangle}} = \frac{\sum_{t_i \in D} \prod_{z_j \in Z} d_{z_j}(c_j, t_i \cdot z_j)}{\sum_{t_i \in D} \prod_{x_j \in X} d_{x_j}(a_j, t_i \cdot x_j)},$$

where $Z = X \cup Y$, $C = A \cup B$.

3 Weighted Quantitative Association Rules

Let itemset $\langle X, A \rangle$ be a pair, where X is the set of attributes x_j and A is the set of fuzzy sets a_j , $j = 1, \dots, p$. We assign a weight $w_{(x,a)}$ for each itemset $\langle X, A \rangle$, with $0 \leq w_{(x,a)} \leq 1$, to show the importance of the item.

Generalizing Definition 2, we can define the weighted fuzzy support for the weighted itemset as follows:

Definition 5. The *weighted fuzzy support* of an itemset $\langle X, A \rangle$ is

$$\begin{aligned} WFS_{\langle X, A, w \rangle} &= \left(\prod_{x_j \in X} w_{(x_j, a_j)} \right) \cdot FS_{\langle X, A \rangle} \\ &= \frac{\sum_{t_i \in D} \prod_{x_j \in X} w_{(x_j, a_j)} d_{x_j}(a_j, t_i \cdot x_j)}{|D|} \end{aligned}$$

Notice the difference from [6] where the sum of weights is used, instead of the product. There are still other ways to define the combined weight. For example, the minimum of item weights would also make sense, and result in a simple algorithm. Product and minimum share the important property that if one item has zero weight, then the whole set has zero weight.

Similar to [8], a support threshold and a confidence threshold will be assigned to measure the strength of the association rules.

Definition 6. An itemset $\langle X, A \rangle$ is called a **frequent itemset** if the weighted fuzzy support of such itemset is greater than or equal to the (user defined) minimum support threshold.

Definition 7. The *weighted fuzzy confidence* of the rule “If X is A then Y is B ” as follows:

$$WFC_{\langle \langle X, A, w \rangle, \langle Y, B, w \rangle \rangle} = \frac{WFS_{\langle Z, C, w \rangle}}{WFS_{\langle X, A, w \rangle}}$$

where $Z = X \cup Y$, $C = A \cup B$.

Definition 8. An fuzzy association rule “If X is A then Y is B ” is called an **interesting rule** if $X \cup Y$ is a frequent itemset and the confidence, defined in definition 7, is greater than or equal to a (user defined) minimum confidence threshold.

The frequent itemsets in the weighted approach have the important property, that all its subsets are also frequent. Thus, we can apply the traditional bottom-up algorithm, tailored to fuzzy sets in [8].

4 Normalized Weighted Quantitative Association Rules

There is one possible problem with our definition. Even if each item has a large weight, the total weight may be very small, when the number of items in an itemset is large.

In this section, we deal with the mining of weighted association rules for which the weight of an itemset is normalized by the size of the itemset.

Definition 9. The *normalized weighted fuzzy support* of itemset $\langle X, A \rangle$ is given by

$$NWFS_{\langle X, A, w \rangle} = \left(\prod_{x_j \in X} w_{(x_j, a_j)} \right)^{1/k} \cdot FS_{\langle X, A, w \rangle}$$

where $k =$ size of the itemset $\langle X, A \rangle$.

Notice that we actually use the geometric mean of item weights as the combined weight. This is a direct analogy to [6], where the arithmetic mean is applied in normalization.

Definition 10. A k -itemset $\langle X, A \rangle$ is called a **frequent itemset** if the normalized weighted fuzzy support of such an itemset is greater than or equal to the minimum support threshold, or

$$NWF_{S_{\langle X, A, w \rangle}} \geq \text{minsup}$$

It is not necessarily true for all subsets of a frequent itemset to be frequent. In [8], we generated frequent itemsets with increasing sizes. However, since the subset of a frequent itemset may not be frequent, we cannot generate candidate k -itemsets simply from the frequent $(k - 1)$ -itemsets. All k -itemsets, which may contribute to be subsets of future frequent itemsets, will be kept in candidate generation process.

We can tackle this problem by using a new z -potential frequent subset for each candidate itemset.

Definition 11. A k -itemset $\langle X, A \rangle$ is called a **z -potential frequent subset** if

$$\left(\prod_{x_j \in X} w_{(x_j, a_j)} \cdot \prod_{y_j \in Y} w_{(y_j, b_j)} \right)^{1/z} \cdot FS_{\langle X, A, w \rangle} \geq \text{minsup}$$

where z is a number between k and the maximum possible size of the frequent itemset, and for each $\langle Y, B \rangle$, $Y \neq X$, is the remaining itemset with maximum weights.

4.1 Algorithm for Mining Normalized Weighted Association Rules

A trivial algorithm would be to solve first the non-weighted problem (weights=1), and then prune the rules that do not satisfy the weighted support and confidence. However, we can take advantage of weights to prune non-frequent itemsets earlier, and thereby cut down the number of trials.

An algorithm for mining normalized weighted quantitative association rules has the following inputs and outputs.

Inputs: A database D , two threshold values minsup and minconf .

Output: A list of interesting rules.

Notations:

D	the database
D_T	the transformed database
w	the itemset weights
F_k	set of frequent k -itemsets (have k items)
C_k	set of candidate k -itemsets (have k items)
I	complete item set
minsup	support threshold
minconf	confidence threshold

```

Main Algorithm ( $minsup, minconf, D$ )
1   $I = Search(D)$ ;
2   $(C_1, D_T, w) = Transform(D, I)$ ;
3   $k = 1$ ;
4   $(C_k, F_k) = Checking(C_k, D_T, minsup)$ ;
5  while ( $|C_k| \neq \emptyset$ ) do
6  begin
7     $inc(k)$ ;
8    if  $k == 2$  then
9       $C_k = Join1(C_{k-1})$ 
10     else  $C_k = Join2(C_{k-1})$ ;
11      $C_k = Prune(C_k)$ ;
12      $(C_k, F_k) = Checking(C_k, D_T, minsup)$ ;
13      $F = F \cup F_k$ ;
14  end
15   $Rules(F, minconf)$ ;

```

The subroutines are outlined as follows:

1. $Search(D)$: The subroutine accepts the database, finds out and returns the complete item set $I = \{i_1, i_2, \dots, i_m\}$.
2. $Transform(D, I)$: This step generates both a new transformed (fuzzy) database D_T from the original database by user specified fuzzy sets and weights for each fuzzy set. At the same time, the *candidate 1-itemsets* C_1 will be generated from the transformed database. If a 1-itemset is *frequent* or *z-potential frequent subset* then it will be kept in C_1 , else it will be pruned.
3. $Checking(C_k, D_T, minsup)$: In this subroutine, the transformed (fuzzy) database is scanned and the weighted fuzzy support of candidates in C_k is counted. If its weighted fuzzy support is larger than or equal to $minsup$, we put it into the frequent itemsets F_k .
4. $Join1(C_{k-1})$: The *candidate 2-itemsets* will be generated from C_1 as in [8].
5. $Join2(C_{k-1})$: This Join step generates C_k from C_{k-1} , similar to [2].
6. $Prune(C_k)$: During the prune step, the itemset will be pruned in either of the following cases:
 - A subset of the candidate itemset in C_k does not exist in C_{k-1} .
 - The itemset cannot be a *z-potential frequent subset* of any frequent itemset.
7. $Rules(F)$: Find the rules from the *frequent itemsets* F as in [2].

5 Experimental Results

We assessed the effectiveness of our approach by experimenting with a real-life dataset. The data had 5 quantitative attributes: monthly-income, credit-limit, current-balance, year-to-date balance, and year-to-date interest.

The experiments will be done applying only the normalized version of the algorithm (as discussed in Section 4), due to reasons explained above. We use the above five quantitative attributes where three fuzzy sets are defined for each of them.

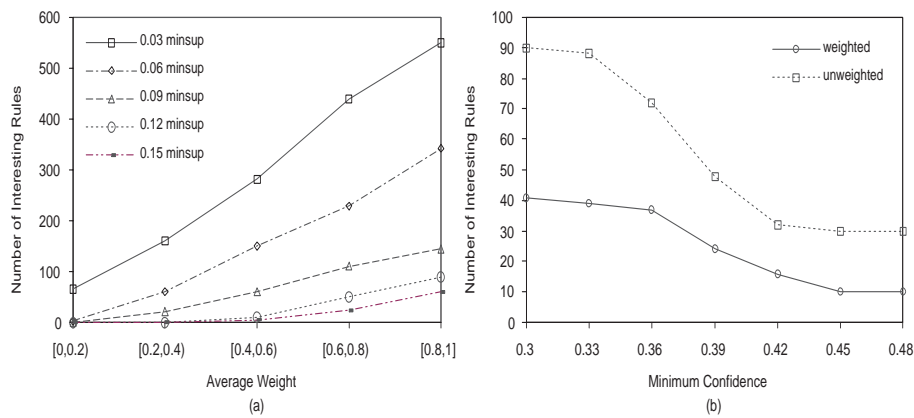


Fig. 1. Left : Random Weight Intervals. Right : Number of Interesting Rules.

Fig. 1(a) shows the increase of the number of rules as a function of average weight, for five different support thresholds. The minimum confidence was set to 0.25. We used five intervals, from which random weights were generated. The increase of the number of rules is close to linear with respect to the average weight. In the following experiments we used a random weight between 0 and 1 for each fuzzy set.

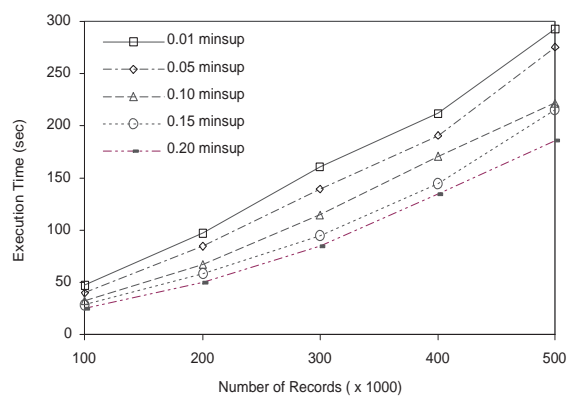


Fig. 2. Scale-up: Number of Records

Fig. 1(b) shows the number of generated rules as a function of minimum confidence threshold, for both weighted and unweighted case (as in [8]). The

minimum support was set to 0.1. The results are as expected: the numbers of rules for the unweighted case are larger, but both decrease with increasing confidence threshold.

Finally, we examined how the performance varies with the number of records. This is confirmed by Fig. 2, which shows the execution time as we increase the number of input records from 100,000 to 500,000, for five different support thresholds. The graphs show that the method scales almost linearly for this dataset.

6 Conclusion

We have proposed generalized mining of weighted quantitative association rules based on fuzzy sets for data items. This is an extension of the fuzzy quantitative association mining problem. In this generalization, the fuzzy sets are assigned weights to reflect their importance to the user. The fuzzy association rule is easily understandable to a human because of the linguistic terms associated with the fuzzy sets.

We proposed two different definitions of weighted support: without normalization, and with normalization. In the normalized case, a subset of a frequent itemset may not be frequent, and we cannot generate candidate k -itemsets simply from the frequent $(k-1)$ -itemsets. We tackled this problem by using the concept of *z-potential frequent subset* for each candidate itemset. We proposed a new algorithm for mining weighted quantitative association rules. The algorithm is applicable to both normalized and unnormalized cases but we prefer the former, as explained above. Therefore, the performance evaluation has been done only for the normalized version.

The results show that by associating weight to fuzzy sets, non-interesting rules can be pruned early and execution time is reduced.

References

1. Piatestky-Shapiro, G., Frawley, W.J.: Knowledge Discovery in Databases. AAAI Press/The MIT Press, Menlo Park, California (1991)
2. Agrawal, R., Imielinski, T., Swami, A.: Mining association rules between sets of items in large databases. Proceedings of ACM SIGMOD (1993) 207–216
3. Agrawal, R., Srikant, R.: Fast algorithms for mining association rules in large databases. Proceedings of the 20th VLDB Conference (1994) 487–499
4. Mannila, H., Toivonen, H., Verkamo, A.I.: Efficient Algorithms for discovering association rules. KDD-94: AAAI Workshop on KDD (1994) 181–192.
5. Toivonen, H.: Sampling large databases for association rules. Proceedings of the 20th VLDB Conference (1996)
6. Cai, C.H., Fu, Ada W.C., Cheng, C.H., Kwong, W.W.: Mining Association Rules with Weighted Items. Proc. of IEEE International DEAS (1998) 68–77.
7. Srikant, R., Agrawal, R.: Mining quantitative association rules in large relation tables. Proceedings of ACM SIGMOD (1996) 1–12.
8. Gyenesei, A.: A Fuzzy Approach for Mining Quantitative Association Rules. Turku Centre for Computer Science, Technical Report No. 336. (2000)