

A method for association rule quality evaluation based on information theory

D. Sitnikov¹, E. Titova¹ & O. Ryabov²

¹*Kharkov State Academy of Culture, Ukraine*

²*National Institute of Advanced Industrial Science and Technology, Japan*

Abstract

The concept of patterns representing functional, logical and other dependencies in data lies in the basis of the Data Mining technology. One of the wide spread forms for representing discovered knowledge patterns is association rules. A method for evaluating an association rule from the viewpoint of information theory has been suggested, which allows us to calculate a generalized characteristic of associations (based on mutual information) with the help of the well known association rule parameters: Support, Confidence and Improvement. Using such a characteristic of associations complements the traditional association parameters and allows setting a linear order on the set of associations, which is useful for evaluating and filtering obtained dependencies. Besides we have carried out analysis of the dependence of the association rule self-descriptiveness on the standard parameters.

1 Introduction

A general definition of association rules has been suggested in [1]: Let $L = I_1, I_2, \dots, I_m$ be a set of object features. Let T be a set of records. Each record t is represented by a binary vector $t[k]=1$ if t contains the feature I_k and $t[k]=0$ if t does not contain the feature I_k ($k = \overline{1, m}$). Let X be a subset including some features from L , i.e. $X \subseteq L$. We say that the record t satisfies X if $\forall I_k \subseteq X, t[k]=1$. An association rule is an expression in the form $X \rightarrow Y$, where $X \subseteq L, Y \subseteq L$, at that $X \cap Y = \emptyset$.



The association rule $X \rightarrow Y$ is supported in the set of records with the confidence level (briefly Conf) c if $c\%$ records in T that contain X also contain Y . The rule $X \rightarrow Y$ has a support (briefly Sup) of s in the set T if $s\%$ records in T contain $X \cup Y$.

To discover association rules a set of algorithms have been developed [1-4, 6, 7], which consist of the major two parts: finding all “large sets” (covers) and forming logic rules by splitting covers into subsets.

There are some problems in evaluating obtained association rules. It has turned out that the well known parameters Support and Confidence are not sufficient to completely characterize obtained association dependencies. Various additional and alternative parameters for evaluating association rules have been suggested in the last years: “interest”, correlation, improvement (quite popular now) [3, 5].

Besides, it seems reasonable to develop a certain integral characteristic that would take into account several parameters. How to compare for example two association rules one of which has greater Support and Improvement and the other has a greater Confidence? Which rule is “better” and how much better? It is interesting to find such a characteristic for the analysis of an association, which firstly would be calculated with the help of the above three parameters and secondly would increase when associations are “stronger” and “decrease” when they are weaker.

1.1 Characteristics of association rules from the probabilistic viewpoint

Support, Confidence and Improvement are generally accepted characteristics of association rules [1-4]. Since association rules are probabilistic by nature it seems natural to consider them from the probabilistic viewpoint. $\text{Supp}(X \rightarrow Y)$ is defined as the ratio of the number of records in the database satisfying the rule $X \rightarrow Y$ to the total number of database records (it can also be defined as the number of records satisfying the given rule). From the probabilistic viewpoint support reflects the probability of the fact that an object (database record) has two features X and Y : $P(XY)$. $\text{Conf}(X \rightarrow Y)$ (rule probability) is defined as the ration of the number of records satisfying the rule to the number of records satisfying the antecedent (left side of the rule). From the probabilistic viewpoint confidence reflects the conditional probability $P_X(Y) = \frac{P(XY)}{P(X)}$. If $X \rightarrow Y$ is an association rule then the following conditions should be met:

$$\text{Supp}(X \rightarrow Y) = P(XY) \geq \min \text{Supp};$$

$$\frac{\text{Supp}(X \rightarrow Y)}{\text{Supp}(X)} = P_X(Y) \geq \min \text{Conf},$$

where $\min\text{Supp}$, $\min\text{Conf}$ are the minimal acceptable values of support and confidence.



Imp (improvement) is defined as follows:

$$\text{Imp}(X \rightarrow Y) = \frac{P_X(Y)}{P(Y)} = \frac{\text{Conf}(X \rightarrow Y)}{\text{Supp}(Y)}.$$

Thus improvement characterizes an increase in the probability of the event Y under the condition that the event X has occurred as compared with the unconditional probability of Y.

Let us demonstrate that the above three parameters allow us to fully characterize (in some sense) a binary association. If both features can take on only values 0 and 1 we can write down the following probabilities P_{00} , P_{01} , P_{10} , P_{11} that give us all possible combinations of feature values and can be expressed in terms of Support, Confidence and Improvement.

First let us express Support, Confidence and Improvement of the rule $X \rightarrow Y$ in terms of the above probabilities:

$$\text{Supp}(X \rightarrow Y) = P(XY) = P_{11};$$

$$\text{Conf}(X \rightarrow Y) = \frac{\text{Supp}(X \rightarrow Y)}{\text{Supp}(X)} = \frac{P(XY)}{P(X)} = \frac{P_{11}}{P_{11} + P_{10}};$$

$$\text{Imp}(X \rightarrow Y) = \frac{\text{Supp}(X \rightarrow Y)}{\text{Supp}(X)\text{Supp}(Y)} = \frac{P(XY)}{P(X)P(Y)} = \frac{P_{11}}{(P_{11} + P_{10})(P_{11} + P_{01})}.$$

From these expressions the probability values can be obtained:

$$P_{11} = \text{Supp}(X \rightarrow Y); \quad (1)$$

$$P_{10} = \frac{\text{Supp}(X \rightarrow Y)}{\text{Conf}(X \rightarrow Y)} - \text{Supp}(X \rightarrow Y) = \text{Supp}(X \rightarrow \bar{Y}); \quad (2)$$

$$P_{01} = \frac{\text{Conf}(X \rightarrow Y)}{\text{Imp}(X \rightarrow Y)} - \text{Supp}(X \rightarrow Y) = \text{Supp}(\bar{X} \rightarrow Y); \quad (3)$$

$$\begin{aligned} P_{00} &= 1 - P_{11} - P_{10} - P_{01} = \\ &= 1 + \text{Supp}(X \rightarrow Y) - \frac{\text{Supp}(X \rightarrow Y)}{\text{Conf}(X \rightarrow Y)} - \frac{\text{Conf}(X \rightarrow Y)}{\text{Imp}(X \rightarrow Y)} = \text{Supp}(\bar{X} \rightarrow \bar{Y}). \end{aligned} \quad (4)$$

Thus in order to calculate P_{00} , P_{01} , P_{10} , P_{11} we have used all the parameters Supp, Conf, Imp and it can be easily seen that if any of these parameters is absent the probability values can not be obtained.

It should be noted that values of Supp, Conf and Imp can not be taken arbitrarily. They should satisfy the following conditions:

$$0 \leq \text{Supp}(X \rightarrow Y) \leq 1. \tag{5}$$

$$\text{Supp}(X \rightarrow Y) \leq \text{Conf}(X \rightarrow Y) \leq 1. \tag{6}$$

$$\frac{(\text{Conf}(X \rightarrow Y))^2}{\text{Supp}(X \rightarrow Y) \cdot \text{Conf}(X \rightarrow Y) + \text{Conf}(X \rightarrow Y) - \text{Supp}(X \rightarrow Y)} \leq \text{Imp}(X \rightarrow Y) \leq \frac{\text{Conf}(X \rightarrow Y)}{\text{Supp}(X \rightarrow Y)}. \tag{7}$$

The above inequalities can be easily proven.

1.2 Evaluating an association from the viewpoint of information theory

We suggest a method for evaluating an association from the viewpoint of information theory. From this point of view the association rule $X \rightarrow Y$ can be considered as information on the event Y obtained as a result of receiving a message about the event X . The information “from event to event” is defined as

$$I_{X \rightarrow Y} = \log_2 \frac{P(XY)}{P(X)P(Y)}, \text{ which is nothing else but } \log_2(\text{Imp}(X \rightarrow Y)).$$

Mutual information is defined as follows:

$$I_{X \leftrightarrow Y} = \sum_{i=1}^n \sum_{j=1}^m P_{ij} \log_2 \frac{P_{ij}}{P_i P_j},$$

where $P_{ij} = P((X \sim x_i)(Y \sim y_j))$ is the probability of the fact that X is in the state x_i and Y is in the state y_j ;

$p_i = P(X \sim x_i)$ is the probability of the fact that X is in the state x_i ;

$p_j = P(Y \sim y_j)$ is the probability of the fact that Y is in the state y_j .

In our case $P_{ij} \in \{P_{00}, P_{01}, P_{10}, P_{11}\}$.

The above considerations allow us to write down a formula for the mutual information of an association, which we call *association self-descriptiveness*:

$$\begin{aligned} I_{X \leftrightarrow Y} = & \text{Supp}(X \rightarrow Y) \cdot \log_2(\text{Imp}(X \rightarrow Y)) \\ & + \text{Supp}(X \rightarrow \bar{Y}) \cdot \log_2(\text{Imp}(X \rightarrow \bar{Y})) + \text{Supp}(\bar{X} \rightarrow Y) \cdot \log_2(\text{Imp}(\bar{X} \rightarrow Y)) \\ & + \text{Supp}(\bar{X} \rightarrow \bar{Y}) \cdot \log_2(\text{Imp}(\bar{X} \rightarrow \bar{Y})). \end{aligned} \tag{8}$$

As we have demonstrated before, the three parameters Sup, Conf and Imp completely describe all the probabilities and therefore all the parameters in the



formula (8) can be expressed in the terms of Sup, Conf and Imp. The expressions for $\text{Supp}(\bar{X} \rightarrow \bar{Y})$, $\text{Supp}(\bar{X} \rightarrow Y)$ and $\text{Supp}(X \rightarrow \bar{Y})$ have been obtained in the formulas (1)-(4). Let us express $\text{Imp}(\bar{X} \rightarrow \bar{Y})$, $\text{Imp}(\bar{X} \rightarrow Y)$ and $\text{Imp}(X \rightarrow \bar{Y})$ in terms of the above association rule characteristics:

$$\text{Imp}(\bar{X} \rightarrow Y) = \frac{\text{Supp}(\bar{X} \rightarrow Y)}{\text{Supp}(\bar{X})\text{Supp}(Y)} = \frac{P_{01}}{(P_{01} + P_{00})(P_{01} + P_{11})}; \quad (9)$$

$$\text{Imp}(X \rightarrow \bar{Y}) = \frac{\text{Supp}(X \rightarrow \bar{Y})}{\text{Supp}(X)\text{Supp}(\bar{Y})} = \frac{P_{10}}{(P_{10} + P_{11})(P_{10} + P_{00})}; \quad (10)$$

$$\text{Imp}(\bar{X} \rightarrow \bar{Y}) = \frac{\text{Supp}(\bar{X} \rightarrow \bar{Y})}{\text{Supp}(\bar{X})\text{Supp}(\bar{Y})} = \frac{P_{00}}{(P_{00} + P_{01})(P_{00} + P_{10})}. \quad (11).$$

Then substitute the expressions for P_{01} , P_{00} , P_{10} , P_{11} from (1)-(4) in (9)-(11) and perform necessary algebraic transformations to get:

$$\text{Imp}(\bar{X} \rightarrow Y) = \frac{\text{Conf}(X \rightarrow Y) - \text{Supp}(X \rightarrow Y) \cdot \text{Imp}(X \rightarrow Y)}{\text{Conf}(X \rightarrow Y) - \text{Supp}(X \rightarrow Y)}; \quad (12)$$

$$\text{Imp}(X \rightarrow \bar{Y}) = \frac{\text{Imp}(X \rightarrow Y) (1 - \text{Conf}(X \rightarrow Y))}{\text{Imp}(X \rightarrow Y) - \text{Conf}(X \rightarrow Y)}; \quad (13)$$

$$\text{Imp}(\bar{X} \rightarrow \bar{Y}) = \frac{1 + \text{Supp}(X \rightarrow Y) - \frac{\text{Supp}(X \rightarrow Y)}{\text{Conf}(X \rightarrow Y)} - \frac{\text{Conf}(X \rightarrow Y)}{\text{Imp}(X \rightarrow Y)}}{\left(1 - \frac{\text{Supp}(X \rightarrow Y)}{\text{Conf}(X \rightarrow Y)}\right) \left(1 - \frac{\text{Conf}(X \rightarrow Y)}{\text{Imp}(X \rightarrow Y)}\right)}. \quad (14)$$

Now we can write down a formula for mutual information, which is expressed in terms of the well known association characteristic: Supp, Conf and Imp.

$$\begin{aligned} I_{X \leftrightarrow Y} &= \text{Supp}(X \rightarrow Y) \cdot \log_2(\text{Imp}(X \rightarrow Y)) \\ &+ \left(\frac{\text{Supp}(X \rightarrow Y)}{\text{Conf}(X \rightarrow Y)} - \text{Supp}(X \rightarrow Y) \right) \cdot \log_2 \left(\frac{\text{Imp}(X \rightarrow Y) (1 - \text{Conf}(X \rightarrow Y))}{\text{Imp}(X \rightarrow Y) - \text{Conf}(X \rightarrow Y)} \right) \\ &\quad + \left(\frac{\text{Conf}(X \rightarrow Y)}{\text{Imp}(X \rightarrow Y)} - \text{Supp}(X \rightarrow Y) \right) \\ &\quad \times \log_2 \left(\frac{\text{Conf}(X \rightarrow Y) - \text{Supp}(X \rightarrow Y) \cdot \text{Imp}(X \rightarrow Y)}{\text{Conf}(X \rightarrow Y) - \text{Supp}(X \rightarrow Y)} \right) \end{aligned}$$

$$\begin{aligned}
 & + \left(1 + \text{Supp}(X \rightarrow Y) - \frac{\text{Supp}(X \rightarrow Y)}{\text{Conf}(X \rightarrow Y)} - \frac{\text{Conf}(X \rightarrow Y)}{\text{Imp}(X \rightarrow Y)} \right) \quad (15) \\
 & \times \log_2 \left(\frac{1 + \text{Supp}(X \rightarrow Y) - \frac{\text{Supp}(X \rightarrow Y)}{\text{Conf}(X \rightarrow Y)} - \frac{\text{Conf}(X \rightarrow Y)}{\text{Imp}(X \rightarrow Y)}}{\left(1 - \frac{\text{Supp}(X \rightarrow Y)}{\text{Conf}(X \rightarrow Y)} \right) \left(1 - \frac{\text{Conf}(X \rightarrow Y)}{\text{Imp}(X \rightarrow Y)} \right)} \right).
 \end{aligned}$$

It should be noted that the expressions in logarithms can not be less than zero, which is due to the inequalities (5-7). If zeros are encountered then the corresponding limit should be used to make calculations $\left(\lim_{P \rightarrow 0} P \log P = 0 \right)$.

It should be noted that the given value (association self-descriptiveness) can be calculated for any probability values, which makes it possible to use it a threshold for filtering association rules (i.e. for cutting off the rules whose self-descriptiveness is less than a given one). If X and Y are independent in the sense that there is no association between them, the rule self-descriptiveness is zero.

The maximum of this parameter $I_{X \leftrightarrow Y} = 1$ is reached if $P_{00}=0,5; P_{01}=0; P_{10}=0; P_{11}=0,5$, i.e. if $\text{Supp}=0,5; \text{Conf} = 1; \text{Imp}=2$. It is obvious that in real world databases it is hardly possible to obtain such parameters. Our tests have shown that normally rules whose self-descriptiveness is greater than 0.2 can be interesting, but the problem of defining numerical thresholds for associations in various types of databases and for different tasks requires much more research.

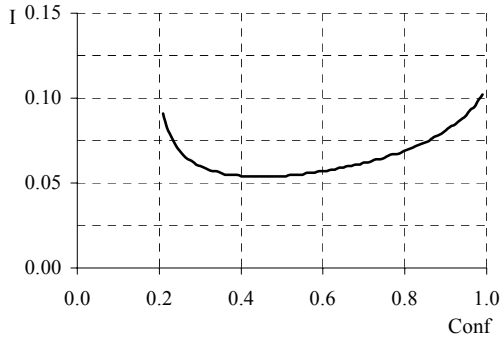


Figure 1.

1.3 Dependence of the association self-descriptiveness on the association rule parameters

Let us consider how Supp, Conf and Imp influence the self-descriptiveness of an association. Traditionally it is known that the greater Supp, Conf and Imp the

better association we have. In order to analyze how the self-descriptiveness of an association depends on the standard characteristics note first that the formula (15) defines a differentiable function with three parameters: Supp, Conf, Imp. Let us draw a graph of this function when Supp and Imp are constant (Supp=0.1; Imp=2), changing Conf in accordance with the inequality (6) (fig.1).

It is obvious that the graph has a minimum. In order to find the value of Conf corresponding to the function minimum the function should be differentiated:

$$\frac{\partial I_{X \leftrightarrow Y}}{\partial Conf} = \frac{Supp}{Conf^2} \times \log_2 \left(\frac{Conf \cdot Imp + Supp \cdot Conf \cdot Imp - Supp \cdot Imp - Conf^2}{Conf \cdot Imp + Supp \cdot Conf \cdot Imp - Supp \cdot Imp - Conf^2 \cdot Imp} \right) + \frac{1}{Imp} \cdot \log_2 \left(\frac{Conf \cdot Imp + Supp \cdot Conf \cdot Imp - Supp \cdot Imp^2 - Conf^2}{Conf \cdot Imp + Supp \cdot Conf \cdot Imp - Supp \cdot Imp - Conf^2} \right).$$

The derivative turns to zero when $Conf = \sqrt{Supp \cdot Imp}$. For example, if Supp=0.1 and Imp=2 the function reaches its minimum if $Conf \approx 0.447$.

Thus the rule with a greater confidence level does not necessarily have a greater self-descriptiveness (i.e. it does not necessarily carry more information from X to Y). For example, when Supp=0.3 and Imp=2 a rule with Conf=0.6 carries more information ($I_{X \leftrightarrow Y} \approx 0.4$) than a rule with Conf=0.78 ($I_{X \leftrightarrow Y} \approx 0.3$), although the confidence level of the latter is 1.3 times higher.

Consider now the influence of Supp on the self-descriptiveness of an association. On differentiating the mutual information function with Conf and Imp being constant we get:

$$\frac{\partial I_{X \leftrightarrow Y}}{\partial Supp} = \frac{1}{Conf} \times \log_2 \left(\frac{Imp(1-Conf)(Conf-Supp)}{Conf \cdot Imp + Supp \cdot Conf \cdot Imp - Supp \cdot Imp - Conf^2} \right) - \log_2 \left(\frac{(1-Conf)(Conf-Supp \cdot Imp)}{Conf \cdot Imp + Supp \cdot Conf \cdot Imp - Supp \cdot Imp - Conf^2} \right).$$

If $Imp > 1$, $\frac{\partial I_{X \leftrightarrow Y}}{\partial Supp} > 0$.

The graph of this function with Conf=0.8; Imp=1.5 illustrates the fact that in this case there are no minimums (fig. 2).

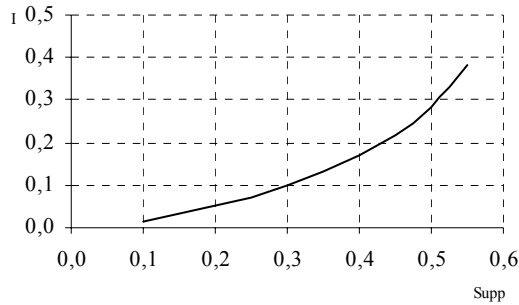


Figure 2.

At last let us investigate how Imp influences the self-descriptiveness of an association.

$$\frac{\partial I_{X \leftrightarrow Y}}{\partial Imp} = \frac{Conf}{Imp^2} \times \log_2 \left(\frac{Conf \cdot Imp + Supp \cdot Conf \cdot Imp - Supp \cdot Imp - Conf^2}{Imp - Conf} \right) - \frac{Conf}{Imp} \cdot \log_2 (Conf - Supp \cdot Imp).$$

This derivative turns to zero when Imp=1, which means that in this case the features are independent in the sense that no information is carried from X to Y (mutual information is zero). An example of the above dependence is shown in fig. 3. Imp changes within the limits defined by the inequality (7).

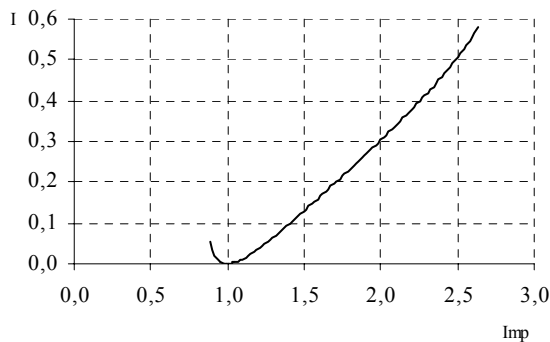


Figure 3.

It should be noted here that Supp and Imp are “symmetric” parameters as far as the “direction” of a rule is concerned, i.e. $\text{Supp}(X \rightarrow Y) = \text{Supp}(Y \rightarrow X)$ and $\text{Imp}(X \rightarrow Y) = \text{Imp}(Y \rightarrow X)$. Therefore when we consider these parameters we can talk about support and confidence of an association and not of an association rule (sometimes these terms are not distinguished). The same is true for the self-descriptiveness of a rule, which is “symmetric” as well. From the viewpoint of information theory it is natural since mutual information from X to Y and from Y to X is always the same. This explains the fact that when Supp and Imp grow the self-descriptiveness grows as well.

Conf is not a “symmetric” parameter, i.e. generally speaking $\text{Conf}(X \rightarrow Y) \neq \text{Conf}(Y \rightarrow X)$. Confidence characterizes the rule itself and the rule “direction” is important.

2 Conclusion

The suggested method for evaluating associations from the viewpoint of information theory is based on calculating an integral characteristic (mutual information), which allows us to compare associations with different values of support, confidence and improvement. Using such a characteristic does not suppose rejecting the use of the standard parameters but complements them and allows taking into account the values of these parameters. The self-descriptiveness of a rule (mutual information) can be used to filter the associations discovered in a database.

The analysis of the dependence of the association self-descriptiveness on the standard characteristics (support, confidence, and improvement) has demonstrated that this dependence is not obvious (it is not always true that the higher the levels of support, confidence and improvement are, the greater the self-descriptiveness is). Thus rules with “bad” values of Sup, Conf or Imp can also be interesting from the viewpoint of the information transferred. Of course, finding associations with “bad” (small) values of support, confidence and improvement is associated with additional expenses (time, memory) and problems associated with the necessity of finding such rules require additional research for particular tasks.

It should be noted that the above considerations are true for binary associations between two features that take on only Boolean values. The cases of more than two Boolean features and arbitrary discrete features are yet to be investigated.

References

- [1] *R.Agrawal, T.Imielinski, A.Swami*. Mining association rules between sets of items in large databases // Proc. of the ACM SIGMOD Conference. – Washington DC, USA, May 1993. – P. 207-216.
- [2] *R.Agrawal, R.Srikant*. Fast algorithms for mining association rules // Proc. of the 20th VLDB Conference Santiago. – Chile, September 1994. – P. 487-499.



- [3] *R.Srikant, R.Agrawal*. Mining generalized association rules // Proc. of the 21st VLDB Conference Zurich. – Switzerland, September 1995. – P. 407-419.
- [4] *A.Amir, R.Feldman, R.Kashi*. A new and versatile method for association generation // Information Systems. – 1997. – Vol. 22, № 6/7. – P. 333-347.
- [5] *Аджиев В.* MineSet – визуальный инструмент аналитика // Открытые системы. – 1997. – №3. – С. 72-77.
- [6] *Hannu Toivonen*. Sampling Large Databases for Association Rules // Proc. of the 22nd International Conference on Very Large Databases. – Mumbai, India, 1996. – P. 134-145.
- [7] *Ashoka Savasere, Edward Omiecinski, Shamkant B. Navathe*. An Efficient Algorithm for Mining Association Rules in Large Databases // Proc. of the 21st International Conference on Very Large Databases. – Zurich, Switzerland, 1995. – P. 432-444.

