

# Multiple uses of frequent sets and condensed representations

## Extended abstract

Heikki Mannila and Hannu Toivonen

University of Helsinki

Department of Computer Science

P.O. Box 26, FIN-00014 University of Helsinki, Finland

Heikki.Mannila@cs.Helsinki.FI, Hannu.Toivonen@cs.Helsinki.FI

### Abstract

In interactive data mining it is advantageous to have condensed representations of data that can be used to efficiently answer different queries. In this paper we show how frequent sets can be used as a condensed representation for answering various types of queries.

Given a table  $r$  with 0/1 values and a threshold  $\sigma$ , a frequent set of  $r$  is a set  $X$  of columns of  $r$  such that at least a fraction  $\sigma$  of the rows of  $r$  have a 1 in all the columns of  $X$ . Finding frequent sets is a first step in finding association rules, and there exists several efficient algorithms for finding the frequent sets. We show that frequent sets have wider applications than just finding association rules.

We show that using the inclusion-exclusion principle one can obtain approximate confidences of arbitrary boolean rules. We derive bounds for the errors in the confidences, and show that information collected during the computation of frequent sets can also be used to provide individual error bounds for each clause. Experiments show that this method enables one to obtain different forms of rules from data extremely fast. Furthermore, we define a general notion of condensed representations, and show that frequent sets, samples and the data cube can be viewed as instantiations of this concept.

### Introduction

Knowledge discovery in databases is an interactive process: one has to be able to look at the data from several different angles. For large data sets, it is not efficient to go back and read the data every time the user wants to see something new from it. A condensed representation of the data for answering most of the needs of the user would be very useful.

One of the most researched areas in data mining is the problem of finding association rules from binary data (Agrawal, Imielinski, & Swami 1993; Houtsma & Swami 1993; Klemettinen *et al.* 1994; Han & Fu 1995;

Holsheimer *et al.* 1995; Park, Chen, & Yu 1995; Savasere, Omiecinski, & Navathe 1995; Srikant & Agrawal 1995; Agrawal *et al.* 1996; Srikant & Agrawal 1996; Toivonen 1996). Assuming we have a relation (table) with 0/1-valued attributes  $R$ , an association rule is an expression  $X \Rightarrow Y$ , where  $X, Y \subseteq R$ . The meaning of such a rule is that whenever a row has a 1 in all the columns of  $X$ , it tends to have a 1 also in all the columns of  $Y$ .<sup>1</sup>

The algorithms for finding association rules all proceed by first finding frequent sets, i.e., sets  $X$  of attributes such that there are sufficiently many rows containing a 1 in each column of  $X$ . From the number of such rows it is easy to compute the confidences of association rules.

In this paper we consider additional uses for frequent sets. We show that the collection of frequent sets can actually serve as a condensed representation of the input data for rules with disjunction and negation, a much larger class of rules than simple association rules. These simple observations lead to a considerable widening in the use of frequent sets. Experiments show that this method enables one to obtain generalized rules from data extremely fast.

We analyze formally the use of frequent sets in answering queries concerning complex rules. Our basic tool is the inclusion-exclusion principle. The frequent sets provide some of the terms of this sum. Using the concept of the border of the collection of frequent sets, we are able to derive bounds on the error caused by omitting the rest of the terms.

Finally, we consider the general theory of condensed representations of data. We define what it means for one class of structures to serve as a representation for another with respect to a class of queries. We dis-

<sup>1</sup>Note that the input data need not be explicitly in 0/1 form: one can just as well consider association rules based on derived attributes. For example, if the relation contains the attribute Age, one can look for rules containing derived binary attributes such as "Age  $\leq$  40".

cuss how the use of frequent sets, samples and the data cube (Gray *et al.* 1996) can be viewed as instances of this general concept. Our concept can be viewed as a generalization of the notion of  $\epsilon$ -approximations from computational geometry (Haussler & Welzl 1987; Mulmuley 1993).

## General rules and frequent sets

In this section we describe a simple and basically well-known rule formalism of general boolean rules, show that knowledge of the frequencies of sets is sufficient to determine rule confidences, and give some examples.

Given a set  $R$  of 0/1-valued attributes  $R$ , a *boolean formula* over  $R$  is a formula  $\varphi$  built from the atomic formulae  $A = 0$  and  $A = 1$ , where  $A \in R$ , using the connectives  $\wedge$ ,  $\vee$ , and  $\neg$ , and parenthesis. Given a 0/1 relation  $r$  with attributes  $R$ , the *frequency*  $s(\varphi, r)$  of  $\varphi$  in  $r$  is

$$\frac{|\{t \in r \mid \varphi \text{ is true of } t\}|}{|r|},$$

i.e., the fraction of rows where  $\varphi$  holds.

We introduce some shorthand notations for having *all* or *at least one* attribute of a set  $X$  equal to 1: For  $X = \{A_1, \dots, A_k\}$  we denote  $a(X, r) = s(A_1 = 1 \wedge \dots \wedge A_k = 1)$  and  $o(X, r) = s(A_1 = 1 \vee \dots \vee A_k = 1)$ ; we simply use notations  $a(X)$  and  $o(X)$  if the relation is clear from the context.

A *boolean rule* over  $R$  is an expression of the form  $\varphi \Rightarrow \psi$ , where  $\varphi$  and  $\psi$  are boolean formulae. The *confidence* of the rule is  $s(\varphi \wedge \psi) / s(\varphi)$ .

The following proposition is of course well-known.

**Proposition 1** For all boolean formulae  $\varphi$  over  $R$  there exists sets  $W_1, \dots, W_m \subseteq R$  and coefficients  $e_1, \dots, e_m \in \{-1, +1\}$  such that for all relations  $r$

$$s(\varphi, r) = \sum_{i=1}^m e_i a(W_i, r).$$

□

That is, knowing the special terms  $a(X, r)$  is sufficient to determine the number of rows satisfying any boolean formula, and hence also the confidence of any boolean rule.

**Example 2** We give some examples of boolean rules and how their confidences can be expressed in terms of the frequencies of conjunctions.

An *association rule* has the form  $X \Rightarrow Y$ , where  $X, Y \subseteq R$ , and the confidence is defined to be  $a(X \cup Y) / a(X)$ .

A rule with *disjunctive right-hand side* has the form  $X \Rightarrow Y \vee Z$  and expresses that if a row  $t \in r$  has a 1 in

each column of  $X$ , then it has a 1 in each column of  $Y$  or in each column of  $Z$ . The confidence of the rule is

$$\frac{a(X \cup Y) + a(X \cup Z) - a(X \cup Y \cup Z)}{a(X)}.$$

Similarly, we can write a rule with a *disjunctive left-hand side*:  $X \vee Y \Rightarrow Z$ . The rule states that if a row  $t \in r$  has a 1 in each column of  $X$  or in each column of  $Y$ , then it has a 1 in each column of  $Z$ . The confidence of the rule is

$$\frac{a(X \cup Z) + a(Y \cup Z) - a(X \cup Y \cup Z)}{a(X) + a(Y) - a(X \cup Y)}.$$

Rules can also have *negation in the left-hand side*, as in  $X \wedge \neg Y \Rightarrow Z$ . This means that if a row  $t \in r$  has a 1 in each column of  $X$  and does not have a 1 in each column of  $Y$ , then it has a 1 in each column of  $Z$ . The confidence of such a rule is

$$\frac{a(X \cup Z) - a(X \cup Y \cup Z)}{a(X) - a(X \cup Y)}.$$

*Negation in the right-hand side* is just as easy. The confidence of a rule of the form  $X \Rightarrow Y \wedge \neg Z$  is

$$\frac{a(X \cup Y) - a(X \cup Y \cup Z)}{a(X)}.$$

□

**Example 3** One of our example data sets  $r$  is an enrollment database of courses in computer science. The data set consists of registration information of 4734 students. There is a row per student, containing the courses the student has registered for. On average, a row has 4 courses. The total number of courses is 127. We first discovered all sets  $X$  of courses with  $a(X, r) \geq 0.05$ . The number of frequent sets is 489. We then computed association rules, rules with disjunctions (two disjuncts) on the left or on the right-hand side, and rules with negations, and then used some simple selection tools to locate interesting rules from these sets.

While already association rules produce lots of interesting and sometimes even surprising information about the data set, looking for more complex rules makes it possible to locate more subtle phenomena. For example, we can notice that the confidences of many rules of the form

$$X \Rightarrow C \text{ programming} \vee C \text{ and Unix}$$

have significantly higher confidences than the individual rules

$$\begin{aligned} X &\Rightarrow C \text{ programming} \\ X &\Rightarrow C \text{ and Unix.} \end{aligned}$$

$a(W, r) \mid W \subseteq R\}$  and the *disjunctive* queries  $Q_v = \{Q'_W : r \mapsto o(W, r) \mid W \subseteq R\}$ .

An  $\varepsilon$ -adequate representation for  $\mathcal{S}$  with respect to  $\mathcal{Q}$  is a class  $\mathcal{R} = \{r_i \mid i \in I\}$  of structures and a query evaluation function  $m : \mathcal{Q} \times \mathcal{R} \rightarrow [0, 1]$  such that for all  $Q \in \mathcal{Q}$  and  $s_i \in \mathcal{S}$  we have

$$|Q(s_i) - m(Q, r_i)| \leq \varepsilon.$$

That is, the values of queries from  $\mathcal{Q}$  on any structure from  $\mathcal{S}$  can be evaluated using the corresponding representation from  $\mathcal{R}$  and the query evaluation function  $m$ .

**Example 10** Consider as the class of structures  $\mathcal{S}_{R,01}$  the class of all 0/1 relations over the set of attributes  $R$ . Consider the query class  $\mathcal{Q}_\wedge = \{Q_W \mid W \subseteq R\}$ , where for  $s \in \mathcal{S}_{R,01}$  we have  $Q_W(s) = a(W, s)$ , i.e., the fraction of rows of  $s$  such that the row has a 1 in each column of  $W$ . Then the collection  $\{Fr(s, \varepsilon) \mid s \in \mathcal{S}_{R,01}\}$  of the frequent sets of  $s$ , for each  $s \in \mathcal{S}_{R,01}$ , provides an  $\varepsilon$ -adequate representation of  $\mathcal{S}_{R,01}$  with respect to  $\mathcal{Q}_\wedge$ .

For the telecommunications alarm database we have a 0.0001-adequate representation that consists only of 128 sets. In the case of our course enrollment database  $r$  the size of a 0.05-adequate representation is  $|Fr(r, 0.05)| = 489$  sets. Note that the number of frequent sets does not depend on the number of rows in the database.  $\square$

The notion of an  $\varepsilon$ -adequate representation is very closely related to the concept of  $\varepsilon$ -approximations and  $\varepsilon$ -nets widely used in computational geometry (Haussler & Welzl 1987; Mulmuley 1993). Roughly, given a set  $N$  and a collection  $\mathcal{S}$  of subsets of  $N$ , an  $\varepsilon$ -approximation is a subset  $M \subseteq N$  such that for each  $X \in \mathcal{S}$  we have

$$\left| \frac{|X \cap M|}{|M|} - \frac{|X|}{|N|} \right| \leq \varepsilon.$$

Thus an  $\varepsilon$ -approximation is a subset that gives a good estimate of the sizes of all subsets in  $\mathcal{S}$ . The  $\varepsilon$ -adequate representations differ in two respects: the representation does not have to be of the same form as the original data, and we have additionally required that such representations exist for every structure in a class of structures. Still, the results on  $\varepsilon$ -approximations help in obtaining results about adequate representations.

The existence result for  $\varepsilon$ -approximations was proved by Vapnik and Chervonenkis, and the sizes of such approximations are naturally connected with the VC-dimension (Vapnik 1982). One can show that samples provide an  $\varepsilon$ -adequate representation for finite query classes.

The data cube (Gray *et al.* 1996) is a recently introduced summarizing representation for relations with arbitrary values. Using the concepts above, the data cube is 0-adequate representation for the class of queries containing all aggregate functions. An interesting possibility is to diminish the space and time complexity of the data cube by allowing some error. It seems that this can be accomplished by using a similar strategy as in the computation of frequent sets.

The preceding discussion of adequate representations is quite tentative: the usefulness of the notion has yet to be conclusively demonstrated. It seems to us, however, that this concept could serve as a unifying point of view to look at several different types of approximate ways of representing information.

## Conclusions

We have shown how the collection of frequent sets can be used as a condensed representation for a relation with 0/1 values. The collection makes it possible to approximate the confidences of arbitrary boolean rules. We have given a strong theorem about the sizes of errors caused by the approximation using the concept of the border of the frequent sets, and have given experimental evidence that the bound is typically extremely good. We have also outlined a possible approach to a general theory of condensed representations and showed how frequent sets, sampling, and also the data cube can be viewed as instances of this concept.

There are several open questions. On the theoretical side, the development of the general notions of condensed representation seems useful. From the practical point of view, a more thorough investigation on the actual sizes of the errors in the approximations could be worthwhile.

## References

- Agrawal, R.; Mannila, H.; Srikant, R.; Toivonen, H.; and Verkamo, A. I. 1996. Fast discovery of association rules. In Fayyad, U. M.; Piatetsky-Shapiro, G.; Smyth, P.; and Uthurusamy, R., eds., *Advances in Knowledge Discovery and Data Mining*. Menlo Park, CA: AAAI Press. 307 – 328.
- Agrawal, R.; Imielinski, T.; and Swami, A. 1993. Mining association rules between sets of items in large databases. In *Proceedings of ACM SIGMOD Conference on Management of Data (SIGMOD'93)*, 207 – 216.
- Grable, D. A. 1993. Sharpened bonferroni inequalities. *Journal on Combinatorial Theory, Series B* 57(1):131 – 137.

- Grable, D. A. 1994. Hypergraphs and sharpened sieve inequalities. *Discrete Mathematics* 132:75 – 82.
- Gray, J.; Bosworth, A.; Layman, A.; and Pirahesh, H. 1996. Data Cube: A relational aggregation operator generalizing group-by, cross-tab, and sub-totals. In *12th International Conference on Data Engineering (ICDE'96)*, 152 – 159.
- Han, J., and Fu, Y. 1995. Discovery of multiple-level association rules from large databases. In *Proceedings of the 21st International Conference on Very Large Data Bases (VLDB'95)*, 420 – 431.
- Hausser, D., and Welzl, E. 1987. Epsilon-nets and simplex range queries. *Discrete Comput. Geom.* 2:127–151.
- Holsheimer, M.; Kersten, M.; Mannila, H.; and Toivonen, H. 1995. A perspective on databases and data mining. In *Proceedings of the First International Conference on Knowledge Discovery and Data Mining (KDD'95)*, 150 – 155.
- Houtsma, M., and Swami, A. 1993. Set-oriented mining of association rules. Research Report RJ 9567, IBM Almaden Research Center, San Jose, California.
- Kahn, J.; Linial, N.; and Samorodnitsky, A. 1995. Inclusion-exclusion: exact and approximate. Manuscript.
- Klemettinen, M.; Mannila, H.; Ronkainen, P.; Toivonen, H.; and Verkamo, A. I. 1994. Finding interesting rules from large sets of discovered association rules. In *Proceedings of the Third International Conference on Information and Knowledge Management (CIKM'94)*, 401 – 407. Gaithersburg, MD: ACM.
- Linial, N., and Nisan, N. 1990. Approximate inclusion-exclusion. *Combinatorica* 10(4):349 – 365.
- Mannila, H., and Toivonen, H. 1996. On an algorithm for finding all interesting sentences. In *Cybernetics and Systems, Volume II, The Thirteenth European Meeting on Cybernetics and Systems Research*, 973 – 978.
- Mulmuley, K. 1993. *Computational Geometry: An Introduction Through Randomized Algorithms*. New York: Prentice Hall.
- Park, J. S.; Chen, M.-S.; and Yu, P. S. 1995. An effective hash-based algorithm for mining association rules. In *Proceedings of ACM SIGMOD Conference on Management of Data (SIGMOD'95)*, 175 – 186.
- Savasere, A.; Omiecinski, E.; and Navathe, S. 1995. An efficient algorithm for mining association rules in large databases. In *Proceedings of the 21st International Conference on Very Large Data Bases (VLDB'95)*, 432 – 444.
- Srikant, R., and Agrawal, R. 1995. Mining generalized association rules. In *Proceedings of the 21st International Conference on Very Large Data Bases (VLDB'95)*, 407 – 419.
- Srikant, R., and Agrawal, R. 1996. Mining quantitative association rules in large relational tables. In *Proceedings of ACM SIGMOD Conference on Management of Data (SIGMOD'96)*.
- Toivonen, H. 1996. Sampling large databases for finding association rules. In *Proceedings of the 22nd International Conference on Very Large Data Bases (VLDB'96)*. To appear.
- Vapnik, V. 1982. *Estimation of Dependencies Based on Empirical Data*. New York: Springer-Verlag.