

# 基于数值属性的关联规则的挖掘

涂星原

(郑州工业大学计自系)

**摘 要** 研究了基于数值属性(Quantitative Attribute)的关联规则的挖掘问题,提出了挖掘算法 Q-Basic 和 Q-AR,试验表明 Q-AR 是有效的。

**关键词** 数据挖掘;关联规则;数值属性

**中图分类号** TP311.13

关联规则(Association Rules)是数据挖掘(Data Mining)研究的重要内容。其应用包括商场购物分析、广告邮寄分析、网络故障分析等。

文献[1~5]对关联规则的挖掘作了有意义的研究。R. Agrawal 等提出了 Apriori 算法<sup>[2]</sup>和挖掘多层次关联规则的 Culmulate, Stratify 等算法<sup>[3]</sup>, J. S. Park 等提出了 DHP 算法<sup>[5]</sup>, J. Han 等提出了面向属性归纳(Attribute-Oriented Induction)的关联规则挖掘算法 ML-T2L1 等<sup>[6]</sup>。这些算法均将事务中项目的数量作为 0/1 值来看待,挖掘布尔关联规则 BAR(Boolean Association Rules)。许多应用中,布尔关联规则不能清楚地表达数据间的关联关系。例如,股票综合指数的上升与价位上升的股票个数的关系。

本文研究了基于数值属性(Quantitative Attribute)的关联规则的挖掘问题和挖掘算法,给出了挖掘基于数值属性的关联规则的算法,基本算法 Q-Basic 和改进算法 Q-AR,并对关联规则的有趣性(Interestingness)作了探讨。实验表明算法 Q-AR 是有效的。

## 1 基于数值属性的关联规则的挖掘问题描述

事务表示为  $\{tid, \langle (A_1, N_1), (A_2, N_2), \dots, (A_j, N_j) \rangle\}$ ,  $tid$  为事务标识号,全局唯一,  $A_i$  是项目 ( $A_i \in I$ ),  $N_i$  是项目  $A_i$  的数量值,  $i=1, 2, \dots, j$ ,  $I$  为项目集合。

定义 1 模式  $P$  定义为  $(A_1, N_1) \cap (A_2, N_2) \dots \cap (A_i, N_i)$ ,  $A_k \in I$ ,  $N_k > 0$  ( $k=1, 2, \dots, i$ )。事务  $t: \{tid, \langle (B_1, M_1), (B_2, M_2), \dots, (B_j, M_j) \rangle\}$ , 若模式  $P$  中任意  $(A_k, N_k)$ , 在  $t$  中有  $(B_s, M_s)$  ( $s=1, 2, \dots, j$ ), 使  $B_s = A_k$ ,  $M_s \geq N_k$ , 称  $t$  包含  $P$ , 或  $t$  支持  $P$ 。

定义 2 模式  $P = (A_1, N_1) \cap (A_2, N_2) \dots \cap (A_i, N_i)$  和  $Q = (A_1, M_1) \cap (A_2, M_2) \dots \cap (A_k, M_k)$ ,  $i < k$ , 对于任意  $j$  ( $j=1, 2, \dots, i$ ):  $N_j < M_j$ , 称为  $P$  小于  $Q$ , 记为  $P < Q$ 。

引理 1 若支持模式  $Q = (A_1, N_1) \cap (A_2, N_2) \dots \cap (A_i, N_i)$ , 模式  $P < Q$ , 则  $t$  支持  $P$ 。

模式  $P$  在事务集合  $D$  中的支持率(support)  $\sigma(P/D) = \frac{D \text{ 中包含模式 } P \text{ 的事务的个数}}{D \text{ 中事务总个数}}$ 。

关联规则  $A \Rightarrow B$  的可信度(confidence)  $\Psi(A \Rightarrow B/D) = \frac{\sigma(AB/D)}{\sigma(A/D)}$ , 其中  $A, B$  均为模式<sup>[1,3]</sup>, 为挖掘有效的关联规则, 定义最小支持率  $\sigma_{\min}$  与最小可信度  $\Psi_{\min}$ 。支持率  $\sigma_{\min}$  的模

收稿日期:1997-09-04;修改稿返回日期:1998-02-15

第一作者 女 1964 年 7 月生 学士学位 讲师

式称为频繁模式, 频繁模式集合记为  $L$ 。包含  $i$  个项目的频繁模式的集合, 称为  $i$  频繁模式集合, 记为  $L_i^{[1]}$ 。

定义 3 若事务  $t$  支持模式  $P=(A_1, N_1) \cap (A_2, N_2) \cdots \cap (A_i, N_i)$ ,  $M_i$  为  $N_i$  与事务  $t$  中  $A_i$  项目的数量值的小者, 称模式  $Q=(A_1, M_1) \cap (A_2, M_2) \cdots \cap (A_i, M_i)$  为模式  $P$  在事务数据中的最大包含模式, 记为  $\text{in-max}(P, t)$ 。若  $t$  不支持模式  $P$ , 模式  $P$  在事务中的最大包含模式  $\text{in-max}(P, t)$  为空。

定义 4 模式集  $S$  的简化集  $\text{design}(S) \{x \mid x \in S, \text{不存在 } y \in S, y < x\}$ 。模式集  $S$  的增长集  $\text{gen}(S) = \{y \mid y \notin S, \text{存在 } x \in S, x < y\}$ 。

定义 5  $i$  简化频繁模式集合  $S_i$  定义为:  $\text{design}(L_i)$ 。

挖掘关联规则即是找出满足最小支持率  $\sigma_{\min}$  与最小可信度  $\Psi_{\min}$  的形如  $A \Rightarrow B$  的规则。挖掘过程分为两步: 首先求出频繁模式集合  $L = U_k L_k$ , 然后由  $L$  构造关联规则集合。R. Agral 给出了由  $L$  求解关联规则的方法<sup>[2]</sup>。本文着重讨论频繁模式的求解。

2 基于数值属性的关联规则的挖掘

下面给出挖掘基于数值属性的关联规则的两个算法: Q-Basic 和 Q-AR。

2.1 基本算法 Q-Basic

算法 Apriori 根据模式  $P$  中的项目是否出现在事务  $t$  中, 判断事务  $t$  是否支持模式  $P$ 。算法 Q-Basic 在模式的表示和事务对模式的支持的判别上, 对算法 Apriori 作如下修改: 模式由包含的项目和项目的数量值组成, 根据定义 1 判断事务  $t$  是否支持模式  $P$ , 即检查模式  $P$  中的项目是否在  $t$  中出现和  $P$  中项目的数量值是否小于  $t$  中相同项目的数量值。

2.2 改进算法 Q-AR

以下假设模式集合中的模式按项目的字典序排列, 相同项目的模式按对应项目的数量值按由小到大排列。模式  $P$  的支持事务的个数记为  $P \cdot \text{support}$ 。

算法 Q-AR(挖掘具有数值属性的关联规则集合):

input: 事务数据库  $D$ , 最小支持率  $\sigma_{\min}$

output: 关联规则集合 AR

method:

$L_1 = 1$  频繁模式集合

$S_1 = \text{degen}(L_1)$

$k = 2$

$L = S_1$

while  $S_{k-1} \not\subseteq \text{NIL}$  Do {

$S\text{-}C_k = \text{simple-apriori-gen}(S_{k-1})$

for every  $t \in D$  do

for every  $Q \in S\text{-}C_k$  do {

$P = \text{in-max}(t, Q)$  //  $P = Q$  与  $t$  的最大包含模式

if  $P \not\subseteq \text{NIL}$  then

if  $P \in S\text{-}C_k$  then  $P \cdot \text{support} = P \cdot \text{support} + 1$

else  $\{S\text{-}C_k = S\text{-}C_k + \{P\}; P \cdot \text{support} = 1\}$

$$L_k = \text{itemset-gen}(S-C_k, \sigma_{\min})$$

$$L = L \cup L_k$$

$$S_k = \text{degent}(L_k)$$

$$k = k + 1\}$$

关联规则集合  $AR = \{A \Rightarrow B \mid A \cap B \in L, \frac{\sigma(A \cap B / D)}{\sigma(A / D)} > \Psi_{\min}\}$

算法 itemset-gen(由候选  $k$  简化频繁模式集合求解  $k$  频繁模式集合  $L_k$ )

input: 候选  $k$  简化频繁模式集合  $S-C_k$

output:  $k$  频繁模式集合  $L_k$

method:

$$L_k = \{\}$$

while  $S-C_k \langle \rangle \text{NIL}$  do {

$P = S-C_k$  中字典序的第一个模式

$P\text{-SET} = S-C_k$  中与  $P$  有相同的项目的模式集合

$S-C_k = S-C_k - P\text{-SET}$

$P\text{-GEN} = \text{gen}(P\text{-SET})$

for every  $Q \in P\text{-GEN}$  do

$Q \cdot \text{support} = 0$

$\text{SET} = P\text{-SET} + P\text{-GEN}$

while  $\text{SET} \langle \rangle \text{NIL}$  do {

$Q = \text{SET}$  中字典序的第一模式

$Q \cdot \text{support} = Q \cdot \text{support} + \sum_{Q \leq P} P \cdot \text{support}$

$\text{SET} = \text{SET} - Q$

if  $Q \cdot \text{Support} \geq \sigma_{\min}$  then  $L_k = L_k + \{Q\}$  }

算法 Simple-apriori-gen 用于求解候选  $k$  简化频繁模式集合。输入:  $k-1$  简化频繁模式集合  $S_{k-1}$ ; 输出  $k$  简化频繁模式集合  $S-C_k$ 。算法 Q-AR 说明: 由  $L_{k-1}$  求  $L_k$  时, 利用 simple-apriori-gen 由  $k-1$  简化频繁模式集合  $S_{k-1}$  构造候选频繁模式集合  $S-C_k$ ; 然后访问事务数据  $D$ , 计算  $S-C_k$  中的候选频繁模式的支持; 利用算法 itemset-gen 由候选  $k-1$  简化频繁模式集合  $S-C_k$  求出  $L_k$ ; 最后求出  $L = \bigcup_k L_k$ 。

算法 itemset-gen 将  $C_k$  中的模式划分为包含相同项目的子集, 对每个子集  $P\text{-SET}$ , 求出其增长集  $P\text{-GEN}$ , 然后根据  $P\text{-SET}$  中模式的支持求出  $P\text{-GEN}$  中模式的支持, 经过  $\min$  的过滤, 最后求出  $k$  频繁模式集合  $L_k$ 。

Q-AR 算法由  $k$  简化频繁模式集合  $S_k$  产生候选频繁模式集合的简化集合  $S-C_k$ , 节约了存储空间, 而且事务  $t$  对模式的支持的判断个数减少, 提高了  $L_k$  的求解速度。

### 3 算法分析

在事务数据个数  $|D| = 14000$ , 项目个数  $|I| = 15$ , 平均项目的数量值  $= 8$ , 事务平均包含的项目数  $= 10$  的试验数据下, 在 Pentium-133 上的试验结果如图 1 和图 2。

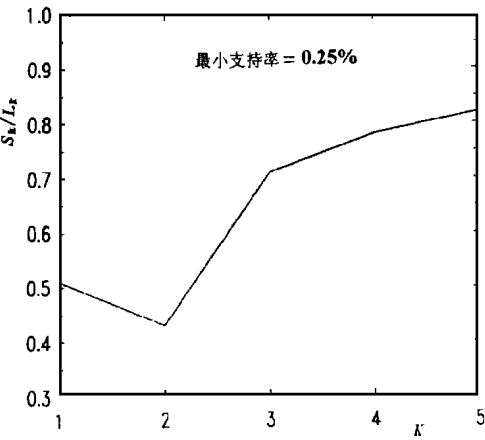


图 1 循环次数  $k$  与  $S_k/L_k$  关系

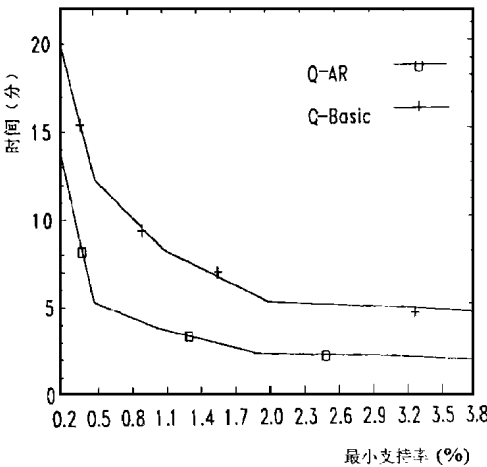


图 2 最小支持率与执行时间关系

4 结论

本文提出了基于数值属性的关联规则的挖掘问题,给出了挖掘算法 Q-Basic 和 Q-AR,并进行了比较,实验表明 Q-AR 是有效的,比 Q-Basic 快 1.5~2 倍,占用存储空间少 1/3。基于数值属性的关联规则是布尔关联规则的扩展,更精确地反映了数据间的关联关系。

参考文献

1 Agrawal, R., Imielinski, T., Iyer, B. et al. Mining association rules between sets of items in larg databases. Proc. of ACM SIGMOD Conf. on management of data, 1993, 207~216

2 Agrawal, R., Srikant, R. Fast algorithms for mining association rules. In the Proc. of VLDB—1994, Chile, Sep. 1994, 487~499

3 Agrawal, R., Srikant, R. mining generalized association rules. In the Proc. of VLDB—1995, Swizerland, 1995, 407~419

4 Fayyad, U., Piatesky-shapiro, G., Smyth, P. From data mining to knowledge discovery: An overview. Advances in knowledge discovery and data mining, MIT Press, 1996, 1~34

5 Park, M., Chen, M., Yu, P. An effective hash based algorithm for mining association rules. in the Proc. of ACM SIGMOD 1995, 175~186

6 Han, J., Fu, Y. Discovery of multiple-level association rules from large databases. in Proc. of VLDB—1995, Swizerland, 1995, 420~431

Efficient Mining of Association Rules Based on Quantitative Attributes

Tu Xingyuan

(Zhengzhou University of Technology)

**Abstract** The mining of the association rules based on quantitative attributes is discussed, and two algorithms, Q-Basic and Q-AR, are proposed in the paper. The experiments show that the algorithm Q-AR is effective.

**Keywords** data mining; association rules; quantitative attribute