

基于十字链表的 Apriori 改进算法

黄建明, 赵文静, 王星星

(西安建筑科技大学信息与控制工程学院, 西安 710055)

摘要: 针对 Apriori 算法中存在的不足, 提出一种把事务数据库映射到十字链表中的改进算法。该算法可以减少连接数据库的次数及事务记录的扫描次数。Apriori 算法与改进算法的性能对比分析表明, 改进算法能有效提高执行效率。

关键词: 数据挖掘; 关联规则; 事务数据库; Apriori 算法; 十字链表

Improved Apriori Algorithm Based on Across Linker

HUANG Jian-ming, ZHAO Wen-jing, WANG Xing-xing

(School of Information and Control Engineering, Xi'an University of Architecture and Technology, Xi'an 710055)

【Abstract】 By studying Apriori algorithm, this paper points out its disadvantages, and proposes the transaction database mapping for an across linker description, which reduces the time cost of the database linking and the scanning time of the transaction record. It gives an improved algorithm. By analyzing the efficiency of the Apriori algorithm and the improved algorithm, it shows that the improved algorithm is better than Apriori algorithm in execution efficiency.

【Key words】 data mining; association rules; transaction database; Apriori algorithm; across linker

1 概述

数据挖掘可以从数据库中发现有价值的、隐藏的知识, 关联规则挖掘^[1]是数据挖掘中一个重要的研究方向, 其主要目标是发现事务数据库中项集之间有趣的关系, 关联规则挖掘过程分为 2 步^[2]: (1)根据一种算法找出所有的频繁项集; (2)根据生成的频繁集产生关联规则。其中, 第(1)步是核心, 决定挖掘关联规则的总体性能, 在确定频繁集后, 就可以得出相应的关联规则, 因此, 算法的研究主要针对第(1)步。

关联规则的相关概念^[3]与性质如下:

定义 1(关联规则) 设 $I = \{I_1, I_2, \dots, I_m\}$ 是所有项的集合, 其中, $I_k (k=1, 2, \dots, m)$ 称为项。项的集合称为项集, 包含 k 个项的项集称为 k -项集。一个事务 T 是一个项集, 它是 I 的一个子集, 每个事务均与一个唯一标识符 Tid 相联系。不同的事务一起组成了事务集 D , 它构成了关联规则发现的事务数据库。关联规则是形如 $X \Rightarrow Y$ 的蕴涵式, 其中, $X, Y \subseteq I$, 且 $X \cap Y = \emptyset$ 。

定义 2(支持度) 关联规则 $X \Rightarrow Y$ 的支持度定义为

$$support(X \Rightarrow Y) = support(X \cup Y) = \frac{|\{T | T \in D \text{ and } (X \cup Y) \subseteq T\}|}{|D|}$$

其中, “ $||$ ”表示集合中的元素个数。

定义 3(频繁项集) 如果项集满足最小支持度, 即如果项集的出现频率大于或等于 min_sup 与事务数据库 D 中事务总数的乘积, 则称该项集为频繁项集, 简称为频繁集, 频繁 k -项集的集合记为 L_k 。

性质 1 频繁项集的所有非空子集也都是频繁的。

2 Apriori 算法

Apriori 算法是最经典的关联规则挖掘算法, 是频繁项集发现的核心算法, 它是一种逐层搜索的迭代方法, 利用已知的频繁 $(k-1)$ -项集 L_{k-1} 生成候选 k -项集 C_k , 再扫描一次数据库来判断候选频繁项集是否是频繁项集, 直到不再产生候选频繁项集, 算法结束。算法描述如下:

输入 事务数据库 D , 最小支持度阈值 min_sup

输出 D 中的频繁项集 L

```
(1)  $L_1 = \text{find\_frequent\_1-itemsets}(D)$ ;  
(2) for  $(k=2; L_{k-1} \neq \emptyset; k++)$  {  
(3)  $C_k = \text{apriori\_gen}(L_{k-1}, min\_sup)$ ;  
(4) for each transaction  $t \in D$  //scan  $D$  for counts  
(5)  $C_t = \text{subset}(C_k, t)$ ; //get the subsets of  $t$  that are candidates  
(6) for each candidate  $c \in C_t$   
(7)  $c.count++$ ;  
(8)  $L_k = \{c \in C_k | c.count \geq min\_sup\}$  }  
(9) return  $L = \bigcup_k L_k$ ;
```

Apriori 算法的第(1)步是找出频繁 1-项集的集合 L_1 , 在第(2)步~第(8)步, L_{k-1} 用于产生候选 C_k , 以找出 L_k 。其中, 第(3)步的 apriori_gen 函数用于产生候选项集, 然后使用 Apriori 性质删除那些具有非频繁子集的候选项集; 产生了所有的候选项集后, 第(4)步扫描数据库; 第(5)步使用 subset 函数找出每个事务中所有的 k -项子集, 并判断这些子集是否与候选项集匹配; 第(6)步、第(7)步则对所有模式匹配成功的候选项集累加计数。

通过对 Apriori 算法的分析可知, 算法存在 2 点不足:

(1) 每生成 L_k 时, 需要扫描数据库一次, 即算法须多次扫描数据库; (2) 在每次扫描数据库的过程中, 候选项集中有些项目或事务不需要扫描, 但算法仍扫描了这些项目和事务。针对上述情况, 本文提出了一种基于十字链表的改进算法。

3 Apriori 算法的改进

为了使算法只扫描一次数据库, 节约存储设备 I/O 时间, 本文将事务数据库中的信息用十字链表表示, 把对数据库的

作者简介: 黄建明(1983 -), 男, 硕士研究生, 主研方向: 数据挖掘, 数据库; 赵文静, 教授; 王星星, 硕士研究生

收稿日期: 2008-07-10 **E-mail:** riluo410@163.com

扫描转变为对内存中十字链表的扫描。

3.1 事务数据库的十字链表表示方法

设事务 $T = \{t_i | 1 \leq i \leq m\}$, 其中 m 为事务 T 的事务长度; $t_i \in I$ 。十字链表的描述可以形式化地表示为: $List = (U, A, V, f)$, 其中 $List$ 为十字链表名; U 为事务的集合形式, $U = \{T_i | 1 \leq i \leq n\}$, T_i 代表某个具体的事务, n 为事务的总数; A 表示每个事务的组合方式, 它由事务 T 与一个唯一标识符 Tid 组成; $V = \{ \langle t_i, number \rangle | t_i \in T, i=1,2,\dots,m, number=m-i+1 \}$, 表示中间节点集合, 其中, 元素 $\langle t_i, number \rangle$ 由事务 T 的项目 t_i 与事务 T 的长度 $number$ 共同组成; f 表示一种映射关系, 定义为 $f: T \rightarrow V$ 。

对于一个经过预处理的数据库模式 $\langle Tid, Item \rangle$, 规定事务中的项按字典序存放, 把此事务数据库映射为十字链表, 则十字链表由事务标识头表、项目头表与链表中中间节点组成, 其中, 每个事务标识节点包含 2 个域: 事务名称和链域 $next$ 指针(指向该事务的第 1 个项); 每个项目节点包含 2 个域: 项目名称和链域 $down$ 指针; 链表中中间节点由 4 个域组成, 如图 1 所示, 包括项目名 $Item$ 、该项目到事务末尾项目的长度 $number$ 、指向下一个事务中此项的 $down$ 指针及指向同一事务中下一项的 $right$ 指针。

$Item$	$number$
$down$	$right$

图 1 十字链表中中间节点表示

例如: 有一事务数据库 D , 其中有 6 条事务记录, 项目名用 a, b, c 等表示, 如表 1 所示。把此事务数据库中的信息用十字链表表示, 如图 2 所示。

表 1 样例数据表

Tid	$Items$
1	a, b, e
2	b, d
3	a, c
4	b, e
5	a, d
6	b, c, d

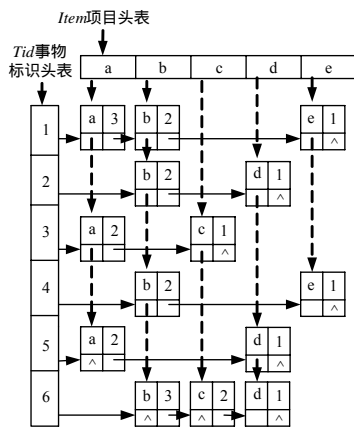


图 2 事务数据库的十字链表表示

3.2 改进的 Apriori 算法

本文基于事务数据库 D 的十字链表的表示, 把在事务数据库中挖掘频繁项集转化为对十字链表的操作, 首先假设: c_i 是候选 k -项集 C_k 中的项集, $c_i[j] (j=0,1,\dots,k-1)$ 表示 c_i 的第 j 项(例如, $c_2[3]$ 表示 c_2 的第 3 项)。

性质 2 若事务 T 中的项与 $c_i[1]$ 匹配成功, 且此项的 $number$ 值小于 k , 则该事务 T 对产生频繁项集是无用的。

很明显, 在生成频繁 k -项集时, 需要扫描十字链表计算

各个候选 k -项集的支持数, 若事务 T 中与 $c_i[1]$ 进行匹配的项的 $number$ 值小于 k , 则事务 T 不可能包含以 $c_i[1]$ 为首的 k -项频繁集的所有项, 因此, 不需要扫描此事务。

性质 3 候选 k -项集中的任一项目集 c_i 与十字链表中的事务匹配时, 假设 $c_i[j]$ 与十字链表中某事务 T 中的项 t 已匹配成功, $c_i[j+1]$ 与 t 的下一项匹配时, 若此项按字典序在 $c_i[j+1]$ 之后, 停止匹配, 此事务对 c_i 的支持数没有影响。

很明显, 如果十字链表中事务 T 中的项已经匹配到字典序在 $c_i[j+1]$ 之后的项, 说明在十字链表的事务 T 中不包含 $c_i[j+1]$, 不能匹配成功, 即对 c_i 的支持数没有影响。

改进算法的步骤如下:

(1) 寻找频繁 1-项集。扫描十字链表, 计算候选 1-项集中各项的支持数; 每项的支持数与给定的 min_sup 进行比较, 如果大于 min_sup , 此项并入频繁 1-项集 L_1 中。

(2) 寻找频繁 k -项集 L_k :

1) 连接 L_{k-1} 产生候选 k -项集, 根据 Apriori 性质对产生的候选 k -项集剪枝后生成候选项集 C_k 。

2) 在候选项集 C_k 中, 沿着十字链表中 $Item$ 项目头表搜索, 将 $c_1[1]$ 与 $Item$ 项目头表中节点的项目名称匹配, 直到匹配成功为止。

3) 沿十字链表中 $Item$ 项目头表中匹配成功节点的 $down$ 指针向下搜索, 判断 $down$ 指针指向的此中间节点的 $number$ 域是否小于 k , 如果 $number$ 值小于 k , 此节点继续沿 $down$ 指针向下搜索, 直到它的 $down$ 指针指向空, 否则, 此节点沿 $right$ 域开始向右匹配 k -项集。

4) 在匹配过程中, 如果项集 c_i 与事务 T 的前 j 项 ($j=0,1,\dots,k-1$) 已经匹配成功, 在进行 $j+1$ 项匹配时, 如果按字典序 $c_i[j+1]$ 位于和事务 T 进行匹配的项之后, 根据性质 3, 停止匹配, 否则, k -项集 c_i 在事务 T 中匹配成功, 说明该 k -项集 c_i 在事务 T 中存在, 则支持数加 1。

5) 判断 k -项集的支持数是否小于 min_sup , 如果不小于, 就将该 k -项集添加到 L_k 中。

6) 取下一个 $c_i[1]$, 进行步骤 2)~步骤 5), 直到 C_k 中的所有 k -项集全部扫描完毕为止。

(3) 输出频繁 k -项集 L_k 。

算法描述如下:

输入 事务数据库 D 映射后的十字链表 $List$, 最小支持度阈值 min_sup

输出 D 中的频繁项集 L

(1) $L_1 = \text{find_frequent_1-itemsets}(List)$;

(2) for ($k=2; L_{k-1} \neq \emptyset; k++$) {

(3) $C_k = \text{apriori_gen}(L_{k-1}, min_sup)$;

(4) for each candidate $c \in C_k$

(5) for each transaction t && firstitem $k_1 \in t$

(6) if (match(c, t)) // 候选 k -项集 c_i 在事务 T 中进行匹配, 匹配成功返回 True

(7) $c.count++$; }

(8) $L_k = \{ c \in C_k | c.count \geq min_sup \}$ }

(9) return $L = \bigcup_k L_k$;

4 算法性能分析

从开始到剪枝, 改进算法与原算法的过程一样。在统计支持数时, Apriori 算法检查每一个事务 T 中是否包含 C_k 中的每一个 k -项集, 需要将 C_k 中的每个 k -项集与事务项 t 匹配,

(下转第 41 页)