

DOI: 10. 3969/j. issn. 1003 - 0972. 2010. 01. 040

一种基于关系代数的 Apriori 优化方法

尤磊*, 兰洋, 熊炎

(信阳师范学院 计算机与信息技术学院, 河南 信阳 464000)

摘要:提出了一种采用关系数据库管理系统的数据库处理能力实现关联规则算法的方法. 结合 Apriori 算法的思想与关系代数的理论, 分析了采用 SQL 语句实现 Apriori 算法的理论可行性, 并描述了算法的实现过程. 在 Mushroom 数据集上的实验验证了本文方法的简单高效性.

关键词:关联规则; 关系代数; SQL 语言; Apriori 算法; 频繁项集

中图分类号: TP311 **文献标志码:** A **文章编号:** 1003-0972(2010)01-0156-05

An Optimized Apriori Implementation Based on Relational Algebra

YOU Lei*, LAN Yang, XIONG Yan

(College of Computer and Information Technology, Xinyang Normal University, Xinyang 464000, China)

Abstract: A method of algorithm implementation for association rules is proposed. It uses the power of data processing in relational database management system to mine data. By combining Apriori algorithm with the theory of relational algebra, the feasibility of implementing Apriori algorithm by use SQL statement is analyzed, and then, the realization process of method is described. An experimental result on mushroom dataset shows that the method given in this paper is simple and efficient.

Key words: association rules; relational algebra; SQL; apriori algorithm; frequent item set

0 引言

数据挖掘是从存放在数据库、数据仓库或其他信息库中的大量数据中挖掘有趣知识的过程^[1]. 挖掘的对象是数据库或数据仓库中的大量数据. 在挖掘过程中, 挖掘程序从数据库或数据仓库中提取数据进行分析.

关系数据库是数据挖掘最流行的、最丰富的数据源, 因此它是我们数据挖掘研究的主要数据形式^[1]. 关系数据库系统的数据库管理系统, 是为数据库的建立、使用和维护而配置的系统软件, 具有强大的数据处理功能^[2].

Apriori 算法是第一个关联规则挖掘算法^[3], 目前对 Apriori 算法的研究主要集中在对算法改进上^[4]. 在实现过程中, Apriori 及改进算法从数据库管理系统中提取数据而忽略了数据库管理系统的的功能. 对于存储在关系数据库或关系数据仓库的数据而言, 可使用关系数据库管理系统的通用

的数据管理功能来实现挖掘算法. 基于此, 本文在分析关联规则的理论及关系代数理论的基础上, 采用关系数据库的 SQL 语言实现 Apriori 算法.

1 相关理论

关联规则 (Association Rules) 挖掘是数据挖掘研究领域的一个重要研究方向, 用于发现隐藏在大型数据集中的令人感兴趣的联系, 发现的联系可以用关联规则或频繁项集的形式表示.

设 $I = \{I_1, I_2, \dots, I_m\}$ 是项的集合, $D = \{T_1, T_2, \dots, T_n\}$ 是一个事务数据库, 其中每个事务 T 是项的结合, 使得 $T \subseteq I$ 每个事务有一个标识符, 称为 TD. 如果 I 的一个子集 X 满足 $X \subseteq T$, 则称事务 T 包含项目集 X . 一个关联规则就是形如 $X \Rightarrow Y$ 的蕴涵式, $X \subseteq I, Y \subseteq I, X \cap Y = \emptyset$.

规则 $X \Rightarrow Y$ 在交易数据库中的支持度 (Support) 就是交易集中包含 X 和 Y 的交易数与所有交

收稿日期: 2009-07-17; 修订日期: 2009-11-04; * 通讯联系人, E-mail: yleiou@163.com

基金项目: 河南省教育厅自然科学基础研究计划项目 (2010A520033); 河南省教育厅青年基金项目 (2009-QN-078)

作者简介: 尤磊 (1981-), 男, 河南罗山人, 硕士, 助教, 主要研究领域为数据挖掘、数据库技术.

易数之比,记为 $s(X \Rightarrow Y)$,即

$$s(X \Rightarrow Y) = \frac{|\{T: X \subseteq T, Y \subseteq T, T \subseteq D\}|}{|D|}$$

规则 $X \Rightarrow Y$ 在交易数据库中的置信度 (Confidence)是指包含 X 和 Y 的交易数与包含 X 的交易数之比,记为 $c(X \Rightarrow Y)$,即 $c(X \Rightarrow Y) = \frac{|\{T: X \subseteq T, Y \subseteq T, T \subseteq D\}|}{|\{T: X \subseteq T, T \subseteq D\}|}$.

关系数据库是以关系模型为基础的数据库.关系数据库的标准语言 SQL 语言是一种介于关系代数和关系演算的语言,其功能包括查询、操纵、定义和控制 4个方面,是一个通用的功能极强的关系数据库标准语言,已经被确定为关系数据库系统的国际标准,被绝大多数商品化的关系数据库系统所采用^[2].

关系模型中的关系操纵能力早期通常是用代数方法或逻辑方法来表示,分别称为关系代数和关系演算.关系代数是通过对关系的运算来表达查询要求的方式,关系演算是用谓词来表达查询要求的方式. SQL 语言是具有关系代数和关系演算双重特点的语言,是一种被关系数据库产品广泛使用的结构化查询语言^[2].

在关系代数中选择运算^[2]是按给定的选择条件选出符合条件的元组,可以表示为:

$$\langle \text{选择条件} \rangle (\langle \text{关系名} \rangle),$$

投影操作选取关系的某些属性,可表示为:

$$\langle \text{属性表} \rangle (\langle \text{关系名} \rangle).$$

根据支持度与置信度的定义以及关系代数中选择运算以及投影运算的定义,可采用关系代数的运算对关联规则的支持度与置信度的描述如下:

$$s(X \Rightarrow Y) = \frac{|\langle \text{同时包含 } X \text{ 和 } Y \text{ 的交易数} \rangle (\langle \text{交易表} \rangle)|}{|\langle \text{交易数} \rangle (\langle \text{交易表} \rangle)|},$$

$$c(X \Rightarrow Y) = \frac{|\langle \text{同时包含 } X \text{ 和 } Y \text{ 的交易数} \rangle (\langle \text{交易表} \rangle)|}{|\langle \text{包含 } X \text{ 的交易数} \rangle (\langle \text{交易表} \rangle)|}.$$

根据以上理论分析,结合关系数据库管理系统强大的数据操纵管理能力,在进行关联分析时,可采用关系数据库语言 SQL 实现关联挖掘算法.

2 Apriori算法

Apriori算法开创性地使用基于支持度的剪枝技术,系统地控制候选项集指数增长,它使用逐层搜索的迭代方法, k 项集用于探索 $(k+1)$ 项集,其主要步骤^[1]有:

(1)连接步:通过 L_{k-1} 与自己连接产生候选 k -

项集的集合,记做 C_k .为方便计,假定事务或项集中的项按字典次序排序,设 l_i 和 l_j 是 L_{k-1} 中的项集. $l_i[j]$ 表示 l_i 的第 j 项, L_{k-1} 与自己连接的条件是它们的前 $k-2$ 个项相同,即

$$(l_i[1]=l_j[1]) \quad (l_i[2]=l_j[2]) \quad \dots \quad (l_i[k-2]=l_j[k-2]) \quad (l_i[k-1] < l_j[k-1]).$$

连接 l_i 和 l_j 的结果项集是 $l_i[1]l_j[2]\dots l_i[k-1]l_j[k-1]$.

(2)剪枝步: C_k 是 L_k 的超集,即 C_k 的成员也可以不是频繁的,但所有的频繁 k 项集都包含在 C_k 中.因一个项集是频繁的,则它的所有子集一定也是频繁的^[3],如果 C_k 的一个候选 k 项集的 $(k-1)$ 子集不在 L_{k-1} 中,则该候选不可能是频繁的,可以从 C_k 中删除,压缩了 C_k .然后扫描数据库确定 C_k 中每个候选项的支持度计数确定频繁项.

Apriori算法的步骤(2)是对步骤(1)产生的候选项集 C_k 中的每一项进行支持度阈值的判断.对压缩后的 C_k ,需要扫描数据库得到支持度计数确定频繁项.若采用常规方法,需要从数据库中读取数据集到内存中,再通过应用程序遍历的方式进行计数.在关系数据库中,可使用关系数据库对数据的操纵管理能力来实现对候选频繁项集的支持度计数.下面结合关系代数的理论以及 Apriori算法的思想进行描述.

设事务数据库的关系模型为 $R(TD, l_1, l_2, \dots, l_m)$,每一个元组的数据采用二元变量^[1]表示.如果项在事务中出现,则它的值为 1,否则为 0,其中 l_1, l_2, \dots, l_m 是项, $I = \{l_1, l_2, \dots, l_m\}$ 是项的集合.根据 Apriori算法步骤(1)中, l_i 和 l_j 是 L_{k-1} 中的项集. $l_i[j]$ 表示 l_i 的第 j 项,则有 $l_i[j] \in I$ 根据关系代数中的选择和投影运算,

$l_i[1]l_j[2]\dots l_i[k-1]l_j[k-1]$ 的支持度计数等价于关系模型 R 中在数据列

$$l_i[1], l_j[2], \dots, l_i[k-1], l_j[k-1]$$

上数据值同时为 1的元组的个数.可设视图 S ,且满足

$$S = \langle l_i[1]=1 \text{ AND } l_j[2]=1 \text{ AND } \dots \text{ AND } l_i[k-1]=1 \text{ AND } l_j[k-1]=1 \rangle (\langle R \rangle),$$

则视图 S 中元组的个数即是 $l_i[1]l_j[2]\dots l_i[k-1]l_j[k-1]$ 的支持度计数.根据关系代数中函数运算^[2]的 COUNT函数,可描述为

$$\text{SELECT COUNT}(\ast) \text{ FROM } S,$$

返回值是 s 中元组的个数,即是

$$l_1[l_1][l_1[2]] \dots l_1[k-1][l_1[k-1]]$$

的支持度计数.

根据以上描述,可综合运用关系代数中的选择运算、投影运算以及函数运算来实现对候选项集的计数.

3 Apriori 算法的 SQL 实现

一个关系模型的逻辑结构就是一张二维表,根据关系的定义,把表 1 中的事务数据库可以看成是一个关系,其关系模型可描述为: $R(TD, \text{面包}, \text{牛奶}, \text{尿布}, \text{啤酒}, \text{鸡蛋}, \text{可乐})$.

在此关系模式中, R 是关系模式的名称,其属性有 TD、面包、牛奶、尿布、啤酒、鸡蛋、可乐. 每一个元组的数据采用表 1 中的二元变量^[1]表示.

表 1 事务数据集的二元表示

Tab 1 Binary transaction data set

TD	面包	牛奶	尿布	啤酒	鸡蛋	可乐
1	1	1	0	0	0	0
2	1	0	1	1	1	0
3	0	1	1	1	0	1
4	1	1	1	1	0	0
5	1	1	1	0	0	1

首先从数据集中挖掘频繁 1 项集,根据关系运算中的函数查询的描述,使用 SUM 函数来统计每项出现的个数,对于 {面包} 项,采用一条 SQL 语句:

```
SELECT SUM(面包) FROM R
```

来实现,返回的结果是关系 R 中“面包”属性在元组值的和,即是项“面包”的计数.

同理,可得出其他几个属性项的计数. 根据 SQL 的语法结构,可实现如下:

```
SELECT SUM(面包), SUM(牛奶), SUM(尿布), SUM(啤酒), SUM(鸡蛋), SUM(可乐) FROM R
```

对表 1 中的事务数据集,可采用如上一条 SQL 语句即可得到每个项目对应的计数. 其返回结果的值为 4, 4, 4, 3, 1, 2 根据支持度计数值 3, 得出频繁 1 项集: {面包, 牛奶, 尿布, 啤酒}.

下一次迭代是根据频繁 1 项集获取频繁 2 项集,根据 Apriori 算法思想,候选频繁 2 项集的个数为 6 个,分别是 {面包, 牛奶}、{面包, 尿布}、{面包, 啤酒}、{牛奶, 尿布}、{牛奶, 啤酒}、{尿布, 啤酒}. 根据关系运算的投影、选择以及函数查询,使

用 COUNT 函数来统计同时包含多个项出现的元组个数:对于 {面包, 牛奶} 项集,实现如下:

```
SELECT COUNT(*) FROM R WHERE
```

```
(面包 = 1) AND (牛奶 = 1),
```

返回结果为

```
3,
```

即 {面包, 牛奶} 项集的计数. 同理有如下 SQL 语句

```
SELECT COUNT(*) FROM R WHERE
```

```
(面包 = 1) AND (尿布 = 1),
```

```
SELECT COUNT(*) FROM R WHERE
```

```
(面包 = 1) AND (啤酒 = 1),
```

```
SELECT COUNT(*) FROM R WHERE
```

```
(牛奶 = 1) AND (尿布 = 1),
```

```
SELECT COUNT(*) FROM R WHERE
```

```
(牛奶 = 1) AND (啤酒 = 1),
```

```
SELECT COUNT(*) FROM R WHERE
```

```
(尿布 = 1) AND (啤酒 = 1).
```

由此得出频繁 2 项集: {面包, 啤酒}, {牛奶, 啤酒}.

根据以上讨论,可得出频繁 k 项集发现频繁 $(k+1)$ 项集的过程. 本文给出 Apriori 算法采用 SQL 实现的算法描述. 设 $R(l_1, l_2, \dots, l_n)$ 为事务数据集的关系模式,其中 l_n 是事务数据集中对应的项, R 的元组数据采用二元变量的形式描述, L_k 为频繁 k 项集, C_k 为候选频繁 k 项集,根据 Apriori 算法的核心思想有 $C_k = L_{k-1} \times L_{k-1}$, $L_k[i][j]$ 表示 L_k 中第 i 个频繁 k 项中的第 j 个项, $len(L_k)$ 表示 L_k 中频繁项的个数. 对于上列中的频繁 2 项集有: $L_2[1][1]$ 表示项面包, $L_2[1][2]$ 表示项啤酒.

算法:采用 SQL 实现 Apriori 算法,使用根据候选生成的逐层迭代得出频繁项集.

输入:事务数据集的关系模式 R ; 最小支持度阈值 min_sup .

输出: R 中的频繁项集 L

方法:

- 1) for ($j = 1; j \leq n; j++$) / 循环对关系 R 中的项计数获取频繁 1 项集
- 2) {
- 3) if (SELECT SUM(l_j) FROM R) $\geq min_sup$ then
- 4) $L_i = L_i \cup \{l_j\}$; / 如果项 l_j 的计数不小于最小支持度则并入频繁 1 项集
- 5) }
- 6) for($k = 2; L_{k-1} \neq \emptyset; k++$) / 根据频繁

```

k - 项集发现频繁 (k + 1) 项集
7) {
8)  $C_k = L_{k-1} \bowtie L_{k-1}$ ; //连接步并压缩  $C_k$ 
9) for( $i = 1$ ;  $i \leq \text{len}(C_k)$ ;  $i++$ ) //循环
    对候选 k项集中的每个项集计数
10) {
11) if (SELECT COUNT(*) FROM R
        WHERE ( $C_k[i][1] = 1$ ) AND ( $C_k[i]$ 
        [ $2$ ] = 1) AND ...AND ( $C_k[i][\text{len}(C_k$ 
        [ $i$ ]) - 1] = 1) AND ( $C_k[i][\text{len}(C_k$ 
        [ $i$ ]) = 1)) >= min_sup then
12)  $L_k = L_k \cup \{C_k[i]\}$ ; //判断  $C_k[i]$ 中的项
    集是否满足支持度计数
13) }
14) }
15) return  $L = \bigcup_k L$ 
    
```

由上面实现过程可知,在采用 SQL 语句实现 Apriori算法时,可根据候选 k项集的项组合形成 SQL 查询语句,提交至数据库管理系统执行后返回计数结果,再由此判断是否属于频繁项.即将确定候选项集计数的任务交给数据库管理系统处理.

4 实验

本文在 Intel酷睿 2双核 T8100处理器,1 024 MB RAM,160 GB 硬盘的联想笔记本上,编程工具为 Delphi7.0,SQLServer2005 数据库、ADO 数据库引擎、采用 <http://archive.ics.uci.edu/ml/datasets/Mushroom> 页面的蘑菇数据库 (Mushroom Data Set) 为测试数据集,展开实验,该数据库有 8 124 条记录,记录了蘑菇的帽子形状、帽子颜色、颈的形状、颈的颜色、气味、生存环境、是否有毒等 23 种属性,每种属性有 2~12 个枚举值.

实验前需将蘑菇数据库预处理为表 1 所示的二元结构表示方式,即采用二进制数据表示记录中某个项出现与否,因蘑菇数据库每种属性有 2~12 个枚举值,需要按照每一个属性 (第一个属性是期望输出类别 (P, E) 除外) 枚举值的个数 m 扩展为 m 个属性,属性的取值代表了项是否出现.预处理后的数据集共有 117 个属性,8 124 条记录,分别采用 SQL 实现的 Apriori 算法与循环数据集对频繁项计数的 Apriori 算法,以随机提取数据集中的 2 000、2 500、3 000 条数据为 3 组测试数据,支持度阈值分别取 40%、50%、60%、70%、80% 进行实验,运行时间 (ms) 比较如图 1 所示.

从图 1 可以看出,在相同记录条数、相同支持度阈值的条件下,采用 SQL 实现的 Apriori 算法执行时间要小于循环数据集计数的方法实现.并且支持度阈值越小,时间相差越大,且从试验数据可以看出,在支持度阈值相同的情况下,记录条数的增加使支持度计数增加,算法执行时间也将有所减少.

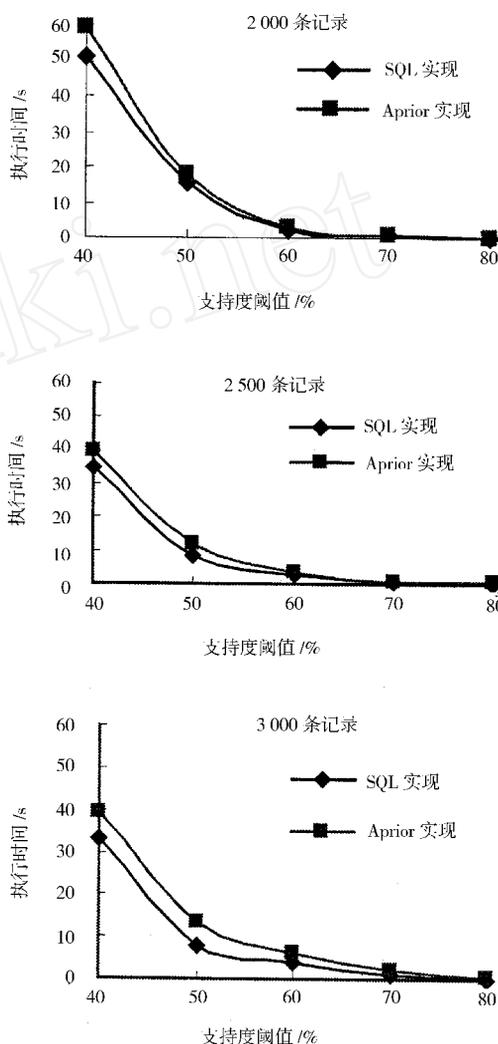


图 1 算法两种实现方法比较

Fig 1 Comparison of two implementation method of algorithm

5 结束语

数据挖掘过程需要处理大量的数据,而处理数据是数据库管理系统的优势所在.在挖掘存储在关系数据库中的数据时,可结合关系数据库系统的数据处理能力采用 SQL 语句来实现挖掘算法.实验

表明,采用 SQL 实现的 Apriori 算法提高了挖掘效率、简化了实现过程。

参考文献:

- [1] Han J W, Kamber M. 数据挖掘概念与技术 [M]. 范明, 孟小峰, 译. 北京:机械工业出版社, 2001.
- [2] 萨师煊, 王珊. 数据库系统概论 [M]. 北京:高等教育出版社, 2000.
- [3] Tan P N, Steinbach M, Kumar V. 数据挖掘导论 [M]. 范明, 范宏建, 译. 北京:人民邮电出版社, 2006.
- [4] Wu X D, Kumar V, Quinlan Ross J, et al. *Top 10 algorithms in data mining* [J]. *Knowledge and Information Systems*(S0219-3116), 2008, 14(1): 1-37.
- [5] Hua K A, Jiang N, Villafane R, et al. *An algebraic approach to system performance analysis using data mining techniques* [C] // *Proceedings of the 2003 ACM symposium on applied computing*, March 09-12, Melbourne, Florida, 2003.
- [6] 文继军, 王珊. SEEKER: 基于关键词的关系数据库信息检索 [J]. *软件学报*, 2005, 16(7): 1270-1281.
- [7] 孟小峰, 周龙骧, 王珊. 数据库技术发展趋势 [J]. *软件学报*, 2004, 15(12): 1822-1836.
- [8] Jeff S. *Mushroom dataset* [EB/OL]. (1987-04-27) [2009-04-18]. <http://archive.ics.uci.edu/ml/datasets/Mushroom>.

责任编辑: 郭红建

(上接第 111 页)

- [5] 丁大均. 关于高层建筑中箱基和上部框架的共同作用 [J]. *建筑结构*, 2002(8): 46-48.
- [6] 罗立平, 赵锡宏. 空间框架结构-厚筏-地基共同作用分析 [C] // *上海高层建筑桩筏与箱基础设计理论*. 上海: 同济大学出版社, 1989.
- [7] 李广信. 有关土的相互作用问题 [J]. *岩土工程学报*, 1996, 18(6): 112-115.
- [8] 谢晨智. 上部结构-筏板基础-地基土相互作用的有限元分析 [D]. 武汉: 湖北工业大学, 2007.

责任编辑: 郭红建

(上接第 123 页)

- [5] Atia A A, Donia A M, Yousif A M. *Synthesis of an ine and thio chelating resins and study of their interaction with Zinc (II), Cadmium (II) and Mercury (II) ions in their aqueous solutions* [J]. *Reactive & Functional Polymers*(S1381-5148), 2003, 56: 75-82.
- [6] 张超灿, 庞金兴, 李曦, 等. 核壳型巯基胺螯合树脂的制备及其吸附性能 [J]. *武汉大学学报: 理学版*, 2001, 47(2): 189-191.
- [7] 李艺. PEI 复合磁性分离材料的制备 [J]. *苏州大学学报: 工科版*, 2003, 23(4): 55-60.
- [8] 忻新泉. 计算机在化学中的应用 [M]. 南京: 南京大学出版社, 1986.
- [9] 张延红, 程国斌, 马伟. 利用 Origin 软件对吸附等温线拟合进行分析 [J]. *计算机与应用化学*, 2005, 22(10): 899-902.
- [10] 程国斌, 张延红, 王康平, 等. 利用 Origin 软件实现化学实验数据的拟合分析 [J]. *化学教学*, 2005(9): 46-48.

责任编辑: 张建安