

相联规则发现的一般性算法研究^{*})

On General Algorithms of Association Rules Discovery

叶阳东 姬安明[✓] 潘玉英 范明
(郑州大学计算机科学系 郑州450052)

摘要 大型事务数据库中相联规则发现是 KDD 中一个很重要的问题。本文描述了相联规则发现的一般性算法,对其核心问题进行了全面和较深入的探讨,并提出了一些提高算法效率的方法。

关键词 相联规则,事务,支持度,置信度,强项集

相联规则发现是数据库应用领域 KDD (Knowledge Discovery in Database) 近期所研究的一个非常重要的问题。在实际的大型事务数据库(如超级市场的事务数据库)中,发现出规则:购买物品 A 和 B 的客户有 95% 的客户购买物品 D 和 C,用规则 $R: A, B \Rightarrow C, D(95\%)$ 来表示。这在分类设计、商店的布局、隶属邮寄、产品的排放、市场分析等多方面大有应用。

相联规则发现算法研究中,一方面借鉴了许多其他发现算法的发现技术,其中包括函数依赖性、强规则、分类规则、因果规则、聚类分析等^[1-3];另一方面,算法在发现多种形式的规则方面进一步完善,这包括一般化规则^[4]、分层规则^[5]等。然而,所有这些发现算法的基本步骤都是相同的,其核心问题在于如何能快速有效地求出强项集。本文主要描述了相联规则发现的一般性算法,并对强项集的求得这一核心问题进行了较深入的探讨,提出用抽样、分区、引入领域知识等方法来提高算法的效率。

一、相关知识

·事务 $I = \{i_1, i_2, \dots, i_m\}$ 是由 m 个不同点组成的集合, i_k 称为项。事务是 I 上的一个子集,也就是由一组点 $i_1, i_2, \dots, i_p \in I$ 构成,每个事务均有一个唯一标识符 Tid,不同事务一起构成了相联规则发现的事务数据库 D 。

·项集 X 的支持度 $\text{Support}(X)$ 表示项集 X 的重要性。设 $X \subset I$ 为项集, $B = D$ 中包含 X 的事务的数量, $A = D$ 中所有事务的数量,项集 X 的支持度

$\text{Support}(X) = B/A$ 。

·最小支持度 发现任务所要求的项集的最小支持度,只有满足于最小支持度的项集才有可能在相联规则中出现,这些项集称之为强项集 (large itemset)。

·规则置信度 (confidence) 规则可靠性的度量。对于规则 $R: X \Rightarrow Y$, 其中 $X, Y \subset I$, 规则 R 的置信度 $\text{confidence}(R) = \text{support}(X \cup Y) / \text{support}(X)$ 。

·相联规则 (association rule) 形式是 $X \Rightarrow Y \{F\}$, 其中 $X, Y \subset I$, 并且 $N \cap Y = \emptyset$, X 称为先决条件, Y 称为规则的结果, F 为规则的置信度。

二、相联规则发现步骤

相联规则发现问题实质上是在满足于最小支持度的项集中,找出满足于最小置信度的相关联的规则。在 D 非常大的情况下,该问题显得尤其复杂,从而成为 KDD 近期研究的热点。所有的发现算法无论它采用什么数据结构,其复杂程度、效率如何,最终都可分为如下几步:

步骤1 预处理与发现任务有关的数据,根据具体问题的要求对数据库进行相应的操作,从而构成规格化后的事务数据库 D 。

步骤2 针对 D , 求出所有满足于最小支持度的项集 $L\text{-set}$ (即强项集)。在 D 很大、 I 中项数 m 非常大的情况下,该问题就构成了发现算法的核心。

步骤3 提取满足于最小置信度的规则集。方法为:

For all itemset $\in L\text{-set}$ and itemset. flag * =

^{*}) 本文得到国家自然科学基金项目和河南省自然科学基金项目部分资助

0and $X \subset \text{itemset}$ and $Y \subset X$ Do

{ $F = \text{support}(X) / \text{support}(X-Y)$; if $F \geq \text{min confidence}$

THEN $Rset = Rset \cup \{X-Y \Rightarrow Y(F)\}$ }

步骤4 解释并输出 Rset。

说明: L.set 中的每一项 itemset 有三个域: set、sup、flag, 分别表示项集、支持度、被包含标志。若 flag=1, 说明在 L.set 中存在有 itemset 的超集; flag=0 就表示 L.set 中无 itemset 的超集。flag 标志域的设置会提高步骤3的速度。

三、算法中的关键性问题--L.set 的产生

由以上的发现步骤可知, 在求取 L.set 时, 由于不同的项集数量可达到 2^n 个, 况且数据库中的事务可能很多, 若对所有的不同项集都进行支持度的计算, 几乎是不可能的。目前有关关联规则地发现算法主要是集中研究如何能快速有效地提取 L.set, 因此其焦点有二, 一是能对项集进行充分地剪支, 尽最大可能地最小化有用项集; 二是以最小次数地扫描事务数据库 D, 从而提高算法的效率。

为了最大可能地减少项集的组合, 所有的发现算法在产生强项集时都遵循这样的原则, 任何强项集的子集都是强项集; 任何弱项集的超集都是弱项集。这样在求强项集时, 算法首先求出项数为一项的强项集 L.1.itemset, 再由 L.1.itemset 产生项数为二的候选项集 C.2.itemset, 扫描 D 计算支持度求出 L.2.itemset, 依次类推产生 C.k.itemset, 扫描 D 求出 L.k.itemset, 其描述如下:

- (1) For all Tid \in D Do {对一项集进行计数求支持度};
- (2) L.1.itemset_i = {支持度大于最小支持度的所有一项集};
- (3) L.set₁ = L.1.itemset;
- (4) k₁ = 2; //k 代表扫描 D 的次数
- (5) While (L.k-1.itemset \neq \emptyset) Do
- (6) {
- (7) C.k.itemset_i = generate.k.itemset from L.k-1.itemset;
- (8) For all Tid \in D Do {对 Tid 所包含的所有 C.k.itemset 中的项集计数求支持度}
- (9)
- (10) L.k.itemset_i = {itemset | itemset \in C.k.itemset and itemset.sup \geq mini-support};
- (11)
- (12) L.set₁ = L.set \cup L.k.itemset;
- (13) K₁ = K + 1
- (14) };
- (15) For all itemset \in L.set and itemset.flag = 0 Do {Forall X \subset itemset and
- (16) X \in L.set Do X.flag = 1} //将 L.set 中的所有 itemset 的子集的 flag 标志置为 1
- (17)

四、提高算法效率的几种方法

由求取 L.set 算法的(7)、(8)、(5)---(14)知, 要

提高算法的效率必须注意以下几个环节:

1. 怎样快速地产生 C.k.itemset。由算法的(6)可知 C.k.itemset 是由不同的 L.k-1.itemset 中的项集结合而产生的。为提高算法的效率可对 L.k-1.itemset 中的项集进行预处理, 对其进行项排序必然会提高 C.k.itemset 的产生速度。

2. 减少扫描事务数据库的次数。由算法的(4)---(13)可知, 若 L.set 中的最大项集的长度是 K(含有 K 个项), 则要对数据库扫描 K 次。怎样减少数据库的扫描次数是提高算法效率的重要途径之一。为减少数据库的扫描次数, 可先对 D 进行预处理, 例如分区, 扫描 D 的分区产生所有的不同长度的候选项, 最后一次性扫描 D, 计数并产生 L.set, 但必须要注意怎样分区, 怎样提高项计数的速度。

3. 怎样提高项集的计数速度。每次扫描数据库都要对 C.k.itemset 中的项集进行测试计数, 可通过对事务的编码和建位图(Bitmap)的方法, 提高每个事务对项集的计数速度。

4. 最大程度地对候选项集进行剪支。该问题是目前许多算法所研究的热点, 在分层关联规则和广义关联规则地发现算法中^[4,5], 利用抽样、预测和分层代码去进行候选项集的优化, 从而使得项集的计数速度和数据库的扫描速度均有较大提高。在算法中引入领域知识去剪除毫无意义的项集是算法提高效率的途径之一。

结束语 出于关联规则的重要发现和发现算法的不完善, 今后一段时间内的许多研究工作仍将围绕着如何提高算法的效率而进行。如何在算法中引入领域知识对算法进行优化, 以及对数据库进行抽样预测优化候选项集, 是我们今后一段时期所工作的目标。

参考文献

- [1] Gregory Pratetsky-Shapiro, Discovery, Analysis, and Presentation of Strong Rules, in G. Pratetsky-Shapiro and W. J. Frawley (eds), Knowledge Discovery in Database, AAAI/MIT Press, 1991
- [2] Jiawei Han et al., Knowledge Discovery in Database: An Attribute-Oriented Approach. In Proc. of the 18th Intl. Conf. on VLDB, Canada, 1992
- [3] Ashok Savasere et al., An Efficient Algorithm for Mining Association Rules in Large Databases, In Proc. of the 21th Intl. Conf. on VLDB, Switzerland, 1995
- [4] Ramakrishnan Srikan, Rakesh Agrawal, Mining Generalized Association Rules, Same to [3]
- [5] Jiawei Han, Yongjian Fu, Discovery of Multiple-Level Association Rules from Large Databases, Same to [3]