

# 基于数据挖掘的课程相关性研究与分析

贾文, 臧明相, 周 鸿

(西安电子科技大学 计算机学院, 陕西 西安 710071)

**摘 要:** 关联规则是数据库中某些特定事件同时发生的概率的简单陈述。关联规则挖掘就是利用特定方法发掘数据库中潜藏的关联规则的过程。文中主要运用数据挖掘中的关联规则和 AprioriTid 算法, 以考务数据库为挖掘对象, 并在挖掘过程中充分运用数据库技术适时地对挖掘数据进行筛选, 有效地提高了挖掘效率。最后, 对课程间的相关性进行了分析和研究, 得到了有效性的结论。

**关键词:** 数据挖掘; 关联规则; AprioriTid 算法

**中图分类号:** TP301.6; G434

**文献标识码:** A

**文章编号:** 1673-629X(2006)12-0178-03

## Research of Curriculum Relevance Based on Data Mining

JIA Wen, ZANG Ming-xiang, ZHOU Hong

(Institute of Computer Science, Xidian University, Xi'an 710071, China)

**Abstract:** Association rule is a simple statement about the cooccurrence probability of certain specific events in the database. Association rule mining is the process which aims to dig the association rule hiding in the database by specific methods. AprioriTid algorithm and association rule in the data mining are addressed. The examination database are employed as the data mining object. The efficiency of data mining are raised by fully used database technique and filtering original data. The relationship among the curriculums are researched, with conclusion satisfied.

**Key words:** data mining; association rule; AprioriTid algorithm

### 1 课程相关性问题的提出

随着科技事业日新月异的进步, 信息技术伴随着互联网技术飞速发展。以信息技术为主的高等院校为适应这一发展, 不断追新地调整着所开设专业的课程, 这种调整不可避免地存在着盲动性和不合理性。教务管理人员以及任课教师却很难直接根据学生的成绩数据分布找出诸如前期课程与后继课程的关系, 并据此进行专业课程的合理设置。因此借助于相应的数据挖掘工具, 发现数据中隐藏的课程相关规律或模式, 并直观地表示相关规律或模式, 为决策提供支持, 是非常必要的。

关联规则是当前数据挖掘研究的主要方式之一。基于教学考务数据库, 运用关联规则挖掘出满足给定支持度和可信度阈值的不同专业、不同课程之间存在的模式或规律, 以指导教务管理人员科学地设置专业课程, 同时也可指导教师合理地安排教学进度及教学内容。例如, 在教学考务数据库中, 可以挖掘出“《高级语言编程》成绩优秀的计算机专业的学生, 《面向对象程序设计》成绩也是优秀的可能性是 63%”, 而电子机械专业学生的这种可能性是 30%。这样, 可以指导教师及学生加强《高级语言编

程》课的教学及学习, 以此来提高《面向对象程序设计》课的教学及学习效果, 甚至是撤消微电子专业的《面向对象程序设计》课程的设置。

### 2 课程相关的关联规则的形式定义

以考务数据库中的学生成绩为例来形式化地描述课程相关的关联规则:

设  $I = \{i_1, i_2, \dots, i_m\}$  是由  $m$  个不同的课程组成的集合,  $D = \{t_1, t_2, \dots, t_n, \dots\}$  为学生的成绩库记录的集合, 其中  $t_i \in I (1 \leq i \leq n)$  表明  $D$  中的每一条记录所包含的课程必在  $I$  中。设  $X, Y$  是一个  $I$  中项的集合;  $T$  是  $D$  中的一个项集 (即满足某种条件的一个集合), 并且  $T \subseteq I$ , 如果  $X \subseteq T, Y \subseteq T^{[1]}$ 。

一个关联规则是形如  $X \Rightarrow Y$  的蕴涵式, 这里  $X \subset I, Y \subset I$ , 并且  $X \cap Y = \emptyset$ 。规则  $X \Rightarrow Y$  在成绩数据库  $D$  中的支持度 (support) 是成绩记录集中包含  $X$  和  $Y$  的记录数与所有记录数之比, 记为  $\text{support}(X \Rightarrow Y)$ , 即

$$\text{support}(X \Rightarrow Y) = |\{T: X \cup Y \subseteq T, T \in D\}| / |D|$$

规则  $X \Rightarrow Y$  在成绩数据库  $D$  中的置信度 (confidence) 是指包含  $X$  和  $Y$  的记录数与包含  $X$  的记录数之比, 记为  $\text{confidence}(X \Rightarrow Y)$ , 即

$$\text{confidence}(X \Rightarrow Y) = |\{T: X \cup Y \subseteq T, T \in D\}| / |\{T: X \subseteq T, T \in D\}|$$

收稿日期: 2006-03-22

作者简介: 贾文 (1967-), 男, 陕西渭南人, 工程师, 研究方向为软件开发、计算机网络。

$$I \mid \{T; X \subseteq T, T \in D\} \mid$$

给定一个学生成绩考务数据库  $D$ , 关联规则挖掘问题就是产生支持度和置信度分别大于用户给定的最小支持度(minsup)和最小置信度(minconf)的关联规则。例如, 在 1 000 条学生成绩记录中有 600 条显示《高级语言编程》成绩优秀, 而这 600 条记录中又有 360 条显示《面向对象程序设计》成绩优秀, 则规则《高级语言编程》成绩优秀就会《面向对象程序设计》成绩优秀的可信度  $C = 360/600 = 0.6$ , 支持度  $S = 360/1000 = 0.36$ 。

从语义的角度来分析, 规则的可信度表示这条规则的正确程度; 支持度表示用这条规则可以推出百分之几的目标, 即这一规则对于整体数据的重要程度。用户可以定义二个阈值, 要求数据挖掘系统所生成的规则的支持度和可信度都不小于给定的阈值。

这样, 就用蕴含式、支持度和可信度唯一标识了每一个挖掘出来的关联规则。例如, 可以这样表示上面提到的例子:

support(《高级语言编程》成绩优秀  $\Rightarrow$  《面向对象程序设计》成绩优秀) = 0.36;

confidence(《高级语言编程》成绩优秀  $\Rightarrow$  《面向对象设计》成绩优秀) = 0.6。

### 3 数据挖掘过程

数据挖掘 (Data Mining) 是数据库知识发现 KDD (Knowledge Discovery in Database) 的核心, 是指从数据库中提取潜在的、有用的、最终可理解的知识的过。关联规则算法则是数据挖掘的一个重要研究方向, 其侧重于确定数据库中不同领域间的联系, 找出满足给定支持度和可信度的多个域之间的相互关系<sup>[2]</sup>, 最具代表性的是经典的 Apriori 算法, 其关联规则的挖掘过程可以分解成以下两个子问题:

(1) 找出成绩数据库  $D$  中所有满足条件的具有用户指定最小支持度的每一个频繁项集  $I$ , 产生所有的  $I$  的非空子集;

(2) 对于  $I$  的每一个非空子集  $s$ ,

若  $\text{support\_count}(I)/\text{support\_count}(s) \geq \text{minconf}$ , 则输出关联规则“ $s \Rightarrow (I - s)$ ”<sup>[3]</sup>。

由于第 2 个子问题较为容易和直观, 目前大量的研究工作主要都集中在第 1 个子问题上。

表 1 为事务数据库  $D$  及其结构。

#### 3.1 挖掘算法

挖掘算法采用经典的 AprioriTid 算法, 考虑到采用纵向数据结构, 在算法上作了一些调整<sup>[4]</sup>。

算法如下:

输入: 事务数据库  $D$ , 最小支持度阈值 minsup。

输出:  $D$  中的频繁项集  $L$ 。

处理流程:

(1)  $\bar{C}_1 = \text{database } D$ ;

表 1 事务数据库  $D$  及其结构

学号 (xh)	姓名 (xm)	性别 (xb)	课程编号 (kcbh)	课程名称 (km)	成绩 (kscj)	专业名称 (zyc)
0301001	习梦丽	女	030701	高数	87	计算机科学
0301002	刘海涛	男	030701	高数	73	计算机科学
0301003	宋文凭	男	030701	高数	81	计算机科学
...	...					
0301001	习梦丽	女	030509	面向对象	92	计算机科学
0301002	刘海涛	男	030509	面向对象	84	计算机科学
0301003	宋文凭	男	030509	面向对象	69	计算机科学
...	...					
0501001	史晓晖	男	050401	电路基础	81	电子机械

(2)  $L_1 = \{C \in \bar{C}_1 \mid C.\text{count}/|D| \geq \text{minsup}\}$ ;  $k = 2$ ;

(3) while( $L_{k-1} \neq \emptyset$ ) do begin //直到不能生成频繁项目集为止

(4)  $C_k = \text{apriori-gen}(L_{k-1})$ ; //生成含  $k$  个元素的候选项目集

(5)  $\bar{C}_k = \emptyset$ ;

(6) For all  $T \in \bar{C}_{k-1}$  do begin

//确定包含在事务中的标志为  $t$ .TID 的  $C_k$  的候选项目集

(7)  $C_t = \{C \in C_k \mid (C - C[k]) \in t.\text{set-of-itemsets} \wedge (C - C[k-1]) \in t.\text{set-of-itemsets}\}$ ;

(8) 计算集合  $C_t$  中增加的所有候选者;

(9) if( $C_t \neq \emptyset$ ) then  $\bar{C}_k += \langle t.\text{TID}, C_t \rangle$ ;

(10) end;

(11)  $L_k = \{C \in C_k \mid C.\text{count}/|D| \geq \text{minsup}\}$ ;

(12) end;

(13) Answer =  $\bigcup_k L_k$ ;

procedure apriori-gen( $L_{k-1}$ )

(1) for each  $p \in L_{k-1}$

(2) for each  $q \in L_{k-1}$

(3) if ( $p.\text{item}_1 = q.\text{item}_1$ )  $\wedge \dots \wedge (p.\text{item}_{k-2} = q.\text{item}_{k-2}) \wedge (p.\text{item}_{k-1} < q.\text{item}_{k-1})$  {

(4)  $c = p \oplus q$ ; //连接两个项集

(5) if has-infrequent-itemset( $c, L_{k-1}$ )

(6) delete  $c$ ; //除去不可能产生频繁项集的候选集

(7) else  $C_k = C_k \cup \{c\}$

(8) }

(9) return  $C_k$

procedure has-infrequent-itemset( $c, L_{k-1}$ )

(1) for each  $(k-1)$ subset  $s$  of  $c$

(2) if  $s \notin L_{k-1}$  return TRUE; else return FALSE

其中,  $D$  表示数据库; minsup 表示给定的最小支持度;

Answer 表示所有频繁项目集。

AprioriTid 方法在由候选频繁项目集确定频繁项目集时只需扫描一遍数据库, 但由于数据库采用纵向结构, 每

