



个人信用数据的关联规则挖掘 教学案例

文件状态:	文件标识:	DMS_03_004
<input type="checkbox"/> 草稿	当前版本:	V01.00.000
<input checked="" type="checkbox"/> 正式发布	作者:	教培部
<input type="checkbox"/> 正在修改	参与者:	研发部
<input type="checkbox"/> 作废	完成日期:	2010-03-18

目 录

1. 概述	3
2. 案例描述.....	3
3. 建模过程.....	3

1. 概述

Apriori 算法是一种最有影响的挖掘布尔关联规则频繁项集的算法。该算法的基本思想是：首先找出所有的频集，这些项集出现的频繁性至少和预定义的最小支持度一样。然后由频集产生强关联规则，这些规则必须满足最小支持度和最小可信度。然后使用第1步找到的频集产生期望的规则，产生只包含集合的项的所有规则，其中每一条规则的右部只有一项，这里采用的是中规则的定义。一旦这些规则被生成，那么只有那些大于用户给定的最小可信度的规则才被留下来。为了生成所有频集，使用了递推的方法。

2. 案例描述

对个人信用数据进行相关性分析。个人信用数据包含下面属性：年龄、性别、地区、收入、婚姻状态、是否生育、是否有车、按揭时是否有抵押、是否有股权计划。利用Apriori 算法挖掘个人信用数据中各个属性之间的关联规则。案例的样本数据如下：

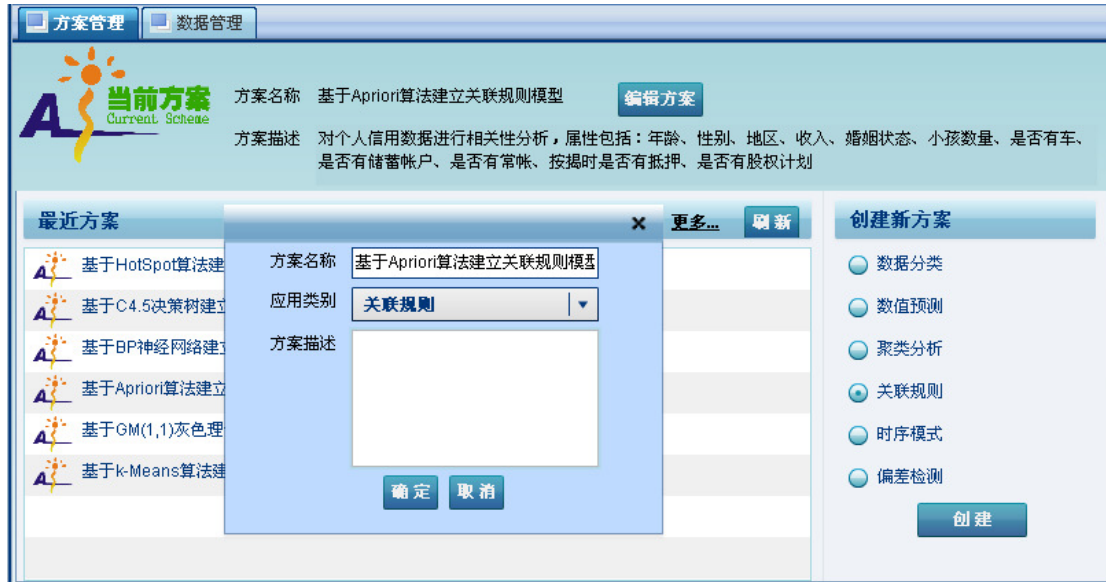
age	sex	region	income	married	children	car	mortgage	pep
48	FEMALE	INNER_CITY	17546	NO	YES	NO	NO	YES
40	MALE	TOWN	30085.1	YES	YES	YES	YES	NO
51	FEMALE	INNER_CITY	16575.4	YES	NO	YES	NO	NO
23	FEMALE	TOWN	20375.4	YES	YES	NO	NO	NO
57	FEMALE	RURAL	50576.3	YES	NO	NO	NO	NO
57	FEMALE	TOWN	37869.6	YES	YES	NO	NO	YES
22	MALE	RURAL	8877.07	NO	NO	NO	NO	YES
58	MALE	TOWN	24946.6	YES	NO	YES	NO	NO
37	FEMALE	SUBURBAN	25304.3	YES	YES	YES	NO	NO
54	MALE	TOWN	24212.1	YES	YES	YES	NO	NO
66	FEMALE	TOWN	59803.9	YES	NO	NO	NO	NO
52	FEMALE	INNER_CITY	26658.8	NO	NO	YES	YES	NO
44	FEMALE	TOWN	15735.8	YES	YES	NO	YES	YES
66	FEMALE	TOWN	55204.7	YES	YES	YES	YES	YES
36	MALE	RURAL	19474.6	YES	NO	NO	YES	NO
38	FEMALE	INNER_CITY	22342.1	YES	NO	YES	YES	NO
37	FEMALE	TOWN	17729.8	YES	YES	NO	YES	NO
46	FEMALE	SUBURBAN	41016	YES	NO	NO	YES	NO
62	FEMALE	INNER_CITY	26909.2	YES	NO	NO	NO	YES
31	MALE	TOWN	22522.8	YES	NO	YES	NO	NO
61	MALE	INNER_CITY	57880.7	YES	YES	NO	NO	YES
50	MALE	TOWN	16497.3	YES	YES	NO	NO	NO
54	MALE	INNER_CITY	38446.6	YES	NO	NO	NO	NO
27	FEMALE	TOWN	15538.8	NO	NO	YES	YES	NO
22	MALE	INNER_CITY	12640.3	NO	YES	YES	NO	NO
56	MALE	INNER_CITY	41034	YES	NO	YES	YES	NO

3. 建模过程

本案例通过太普数据挖掘套件 (<http://www.tipdm.cn>) 实现建模过程。

更多关于此软件工具的介绍详见: <http://www.tipdm.com>

◇ 方案管理



数据管理



预测建模

1、选择算法

从菜单中选择指定算法:

3、参数设置

The screenshot shows the 'Apriori关联规则' (Apriori Association Rules) window. At the top, there are tabs for '方案管理' (Case Management), '数据管理' (Data Management), and 'Apriori关联规则' (Apriori Association Rules). Below the tabs, the window title is '关联规则(Apriori)'. The main area displays a table of training data with columns: sex, region, married, children, and car. A dialog box for parameter settings is overlaid on the table, containing the following fields:

列索引	-1
增量	0.05
最小置信度	0.9
规则条数	10
显著性水平	-1.0
最小支持度下界	0.1
最小支持度上界	1.0

At the bottom of the window, there are buttons for '导入数据' (Import Data), '参数设置' (Parameter Settings), and '规则挖掘' (Rule Mining). The status bar at the bottom indicates the current location: '当前位置 | Apriori关联规则'.

1、规则挖掘

对导入的样本数据使用Apriori关联规则算法对数据进行挖掘。

关联规则(Apriori)

训练数据						
sex	region	married	children	car	mortgage	pep
FEMALE	INNER_CIT	NO	YES	NO	NO	YES
MALE	TOWN	YES	YES	YES	YES	NO
FEMALE	INNER_CIT	YES	NO	YES	NO	NO
FEMALE	TOWN	YES	YES	NO	NO	NO
FEMALE	RURAL	YES	NO	NO	NO	NO
FEMALE	TOWN	YES	YES	NO	NO	YES
MALE	RURAL	NO	NO	NO	NO	YES
MALE	TOWN	YES	NO	YES	NO	NO
FEMALE	SUBURBAN	YES	YES	YES	NO	NO
MALE	TOWN	YES	YES	YES	NO	NO
FEMALE	TOWN	YES	NO	NO	NO	NO
FEMALE	INNER_CIT	NO	NO	YES	YES	NO
FEMALE	TOWN	YES	YES	NO	YES	YES
FEMALE	TOWN	YES	YES	YES	YES	YES
MALE	RURAL	YES	NO	NO	YES	NO
FEMALE	INNER_CIT	YES	NO	YES	YES	NO
FEMALE	TOWN	YES	YES	NO	YES	NO

Apriori
=====

Minimum support: 0.1 (30 instances)
Minimum metric <confidence>: 0.9
Number of cycles performed: 18

Generated sets of large itemsets:

Size of set of large itemsets L(1): 15
Size of set of large itemsets L(2): 86
Size of set of large itemsets L(3): 144
Size of set of large itemsets L(4): 43

Best rules found:

1. children=NO mortgage=NO pep=NO 49 ==> married=YES 48 conf.(0.98)

=== Evaluation ===

Elapsed time: 0.172s

当前位置: Apriori关联规则

更多关联规则挖掘应用参见：<http://www.tipdm.com/Info.php?barbarism=2>