

# 基于时间序列模型的动态关联规则元规则挖掘

刘俊, 张忠林, 谢彦峰, 米伟

(兰州交通大学电子与信息工程学院, 兰州 730070)

**摘要:** 针对现有关联规则挖掘算法大多是挖掘一种静态关联规则的情况, 介绍动态关联规则的定义, 给出动态关联规则元规则的形式化定义, 解决规则随时间的推移可能会有很大变化的情况下为规则建立元规则的问题, 描述一种基于时间序列模型的预测和分析动态关联规则的元规则的方法, 从而较好地拟合历史数据, 给出满足一定显著性水平预测趋势模型的方程, 挖掘规则的变化趋势, 为规则建立元规则。  
**关键词:** 动态关联规则; 元规则; 时间序列; 预测

## Meta-association Rule Mining for Dynamic Association Rule Based on Time Series Model

LIU Jun, ZHANG Zhong-lin, XIE Yan-feng, MI Wei

(School of Electronic and Information Engineering, Lanzhou Jiaotong University, Lanzhou 730070)

**【Abstract】** Almost all of association rules mining algorithms are considered as static ones. In fact, it is possible that a rule will change greatly with time, this paper introduces a rule called dynamic association rule and the definition of the meta-association rules for dynamic association rule, focuses on introducing the method which by the model of time series to mine the meta-association rules for dynamic association rule, and proves that this method is fit on the historical data. It can establish the equation to forecast the tendency of changing rules and mine the meta-association rules for dynamic association rule.

**【Key words】** dynamic association rule; meta-association rule; time series; forecast

为了描述关联规则随时间变化的性质, 文献[1]提出了数据库中动态关联规则新技术, 文献[2]在分析文献[1]原有动态关联规则定义的不足的基础上, 提出了一种新的动态关联规则, 该定义的支持度向量, 置信度向量与经典定义相吻合, 可以更好地反映规则随时间变化的动态信息, 本文在此基础上给出了动态关联规则元规则的形式化定义, 并利用时间序列模型等综合方法对动态规则的元规则进行了分析和预测。

### 1 相关定义和定理

#### 1.1 动态关联规则定义

动态关联规则<sup>[1-2]</sup>是一种能够描述自身特性随时间变化的关联规则, 描述如下:

设  $I = \{i_1, i_2, \dots, i_m\}$  是项集合, 任务相关的事务数据集  $D$  是在时间段  $t$  内收集到的, 同时根据  $t$  的划分, 整个数据集  $D$  可以被分为  $n$  个数据子集:  $D = \{D_1, D_2, \dots, D_n\}$ , 其中, 数据子集  $D_i (i \in \{1, 2, \dots, n\})$  的数据是在  $t_i (i \in \{1, 2, \dots, n\})$  时段内收集的。项集  $T$  满足  $T \subseteq I$ 。若  $A$  和  $B$  为项集,  $A \subset I, B \subset I$ , 并且  $A \cap B = \emptyset$ , 则有如下动态关联规则相关定义。

支持度向量(SV): 动态关联规则  $A \Rightarrow B$  (或者项集  $A \cup B$ ) 的支持度向量具有如下的表示形式:

$$SV = [s_{(A \cup B)_1}, s_{(A \cup B)_2}, \dots, s_{(A \cup B)_n}]$$
$$s_{(A \cup B)_i} = f_{(A \cup B)_i} / |D_i| (i \in \{1, 2, \dots, n\}) \quad (1)$$

其中,  $f_{(A \cup B)_i}$  为项集  $A \cup B$  在数据子集  $D_i (i \in \{1, 2, \dots, n\})$  中出现

的频数;  $|D_i|$  为  $D_i$  中的事务数。设项集  $A \cup B$  的支持度为  $s$ , 则有

$$s = s_{(A \cup B)} = f_{(A \cup B)} / M = \sum_{i=1}^n f_{(A \cup B)_i} / M \quad (2)$$

其中,  $M$  是  $D$  中的事务数。有时项集出现的频数表示支持度更为合适, 这样项集的支持度向量为

$$SV = [f_1, f_2, \dots, f_n]$$

置信度向量(CV): 动态关联规则  $A \Rightarrow B$  的置信度向量具有如下的表示形式:

$$CV = [c_{(A \cup B)_1}, c_{(A \cup B)_2}, \dots, c_{(A \cup B)_n}]$$
$$\text{s.t. } c_{(A \cup B)_i} = s_{(A \cup B)_i} / s_{A_i} \quad (3)$$

其中,  $s_{(A \cup B)_i}$  为项集  $A \cup B$  的 SV 中的第  $i$  个元素,  $s_{A_i}$  为项集  $A$  的 SV 中的第  $i$  个元素。

这样一条完整的动态关联规则就可以表述为: 具有支持度向量 SV、置信度向量 CV、支持度  $s$ 、置信度  $c$  4 个参数的关联规则。它具有如下表示形式:  $A \Rightarrow B(SV, CV, s, c)$ , 其中, SV, CV,  $s, c$  可以分别根据式(1)~式(3)计算, 并一起描述关联规则的动态性质。

#### 1.2 元规则的定义

一般, 元规则形成一个用户希望探察或证实的、感兴趣

**作者简介:** 刘俊(1985-), 男, 硕士研究生, 主研方向: 数据挖掘; 张忠林, 教授、博士; 谢彦峰, 高级工程师; 米伟, 硕士研究生  
**收稿日期:** 2008-12-20 **E-mail:** liujun\_1024@163.com

联系的假定。元规则<sup>[3]</sup>是形如： $P_1 \wedge P_2 \wedge \dots \wedge P_l \Rightarrow Q_1 \wedge Q_2 \wedge \dots \wedge Q_r$ 的规则模板。其中， $P_i (i=1, 2, \dots, l)$ 和 $Q_j (j=1, 2, \dots, r)$ 是示例谓词或谓词变量。

动态关联规则元规则的形式化定义如下：

**定义** 在数据集  $D = \{D_1, D_2, \dots, D_n\} (i \in \{1, 2, \dots, n\})$  上，规则  $A \Rightarrow B$  的支持度向量元规则定义为： $SV = [s_{(A \cup B)_1}, s_{(A \cup B)_2}, \dots, s_{(A \cup B)_n}]$  或  $SV = [f_1, f_2, \dots, f_n]$ ，当且仅当，存在  $i (1 \leq i \leq n)$ ，使  $s_{(A \cup B)_i} \geq \min\_sup$ ，其中， $\min\_sup$  为给定的最小支持度，则此动态规则在数据集  $D$  上的支持度元规则为

$$A \Rightarrow B : SV$$

相似地，在数据集  $D = \{D_1, D_2, \dots, D_n\} (i \in \{1, 2, \dots, n\})$  上的置信度元规则为

$$A \Rightarrow B : CV$$

元规则是“规则的规则”，可以指导知识的发现，帮助提高挖掘过程的性能。元规则可以根据分析者的经验、期望或对数据的直觉或者数据库模式自动生成。

## 2 动态关联规则元规则挖掘描述

现有的关联规则挖掘算法多为静态挖掘，它们的前提假设是：数据集中各项具有近似的性质和作用，即重要性相同或相近；数据集中各项的分布是均匀的，即出现的频率相近或相似<sup>[4]</sup>。然而，不同项的价值和重要程度是不一样的，不同的项在数据集中出现的频繁程度也不一样，因此，利用统一的最小支持度进行整个数据集不同时间段的关联规则挖掘是不充分的。动态关联规则能够提供规则与时间相关的信息，这样便能在一定程度上解决上述不足之处。但动态关联规则仍然存在如何选取支持度的问题，而元规则挖掘的目的之一是发现关联规则序列随时间变化的趋势，通过对趋势的分析，预测下一个时间段的支持度和置信度的可能值。因此，为动态关联规则建立了元规则，在元规则的指导下，可以对数据集中的关联模式进行更加准确有效的挖掘。

目前建立元规则的方法主要有基于概率统计的方法和基于模糊决策树的方法<sup>[5]</sup>。基于概率的方法主要采用回归分析等对规则的支持度进行曲线拟合，这在处理不确定数据上效果欠佳；而从批量处理规则的角度出发，由于需要较多的专家信息，因此基于模糊决策树的方法明显无法满足要求。在历史数据较少时，可以采用基于灰色模型的方法，对于有较多历史数据的情况，本文提出一种基于时间序列模型的方法，可以很好地拟合历史数据，给出满足一定的显著性水平下计算趋势模型的方程。

### 2.1 时间序列模型基本组成

将支持度向量视为一个时间序列，利用时间序列方法进行建模分析并预测。支持度向量的时间序列模型可用叠加形式表示：

$$\hat{S}_t = \hat{f}_t + \hat{p}_t + \hat{x}_t \quad (4)$$

其中， $\hat{S}_t$  为支持度向量中元素  $S_t$  的估计值； $\hat{f}_t$  为趋势分量  $f_t$  的估计值； $\hat{p}_t$  为周期分量  $p_t$  的估计值； $\hat{x}_t$  为随机分量  $x_t$  的估计值。

建立模型的过程就是从已知序列  $S_t (t=1, 2, \dots, n)$  中提取各分量的过程，提取的顺序为趋势分量、周期分量和随即分量。在建立各分量的数学模型后，将其线性叠加，就得到式(4)形式的元规则预测模型。

## 2.2 时间序列模型各分量的确定

### 2.2.1 趋势分量的确定

对于趋势分量  $\hat{f}_t$  可用多项式逼近，即

$$\hat{f}_t = c_0 + c_1 t + c_2 t^2 + c_3 t^3 + \dots + c_k t^k = \sum_{k=0}^k c_k t^k \quad (5)$$

对式(5)可采用多元回归的方法确定系数  $c_0, c_1, c_2, \dots, c_k$  和阶数  $k$ 。具体求解方法可以借用 Excel 软件中县城的回归分析模板来实现。为检验拟合结果，需计算趋势曲线拟合的相关系数  $R$ ，即

$$R^2 = 1 - \frac{\sum_{t=1}^n (S(t) - \hat{S}(t))^2}{\sum_{t=1}^n (S(t) - \bar{S})^2}$$

其中， $n$  为历史数据序列  $S(t)$  的总个数； $\bar{S}$  为序列  $S(t)$  的均值。 $R$  越接近 1， $\hat{S}(t)$  与  $S_t (t=1, 2, \dots, k)$  的线性关系越密切。对给定可信度  $a$  及不同的自由度，可以求出  $R$  的临界值(查表)，只有当  $R$  值大于相应临界值时，回归方程才有实用意义。

### 2.2.2 周期分量的确定

趋势函数确定后，再对扣除趋势分量后的部分  $y_t$  进行周期项分析，即

$$y_t = S_t - \hat{f}_t, t=1, 2, \dots, n$$

本文主要采用谐波分析方法进行周期分量的分析提取。对序列  $y_t$  用  $L$  个波叠加的形式表示其估计值：

$$\hat{p}_t = \frac{a_0}{2} + \sum_{k=1}^L (a_k \cos \frac{2\pi kt}{n} + b_k \sin \frac{2\pi kt}{n}) \quad (6)$$

其中， $\hat{p}_t$  是序列  $y_t$  的估计值，将其作为式(4)的周期分量的估计值； $L$  为谐波的个数，一般取  $n/2$  的整数部分； $k$  为谐波序号，一般认为  $k$  个分波各有  $n/1, n/2, \dots, n/k$  的周期，即第  $k$  个分波的频率为  $k/n$ ； $a_k, b_k$  为傅里叶系数，其计算公式为

$$a_k = \frac{2}{n} y_t \cos \frac{2\pi kt}{n}, b_k = \frac{2}{n} y_t \sin \frac{2\pi kt}{n}, k=0, 1, \dots, L \quad (7)$$

为节省工作量，通常在  $L$  个波中选取波动比较明显的几个谐波相加来估计，在实际应用中只需选取前 6 个显著谐波就能满足精度要求了。若  $s_k^2 = a_k^2 + b_k^2 > 4s^2 \frac{\ln k}{n}$ ，则认为第  $k$  个波显著，否则不显著。其中， $s^2$  为系列的方差，其计算式为  $s^2 = \frac{1}{n-1} \sum_{t=1}^n (x_t - \bar{x}_t)^2$ ； $a$  为检验的显著性水平(一般取 5%)。

### 2.2.3 随机分量的确定

在消除趋势分量和周期分量后，得到随机分量的时间序列，设为  $x_t$ ，即

$$x_t = S_t - \hat{f}_t - \hat{p}_t$$

$x_t$  一般为随机序列，可使用自回归模型(AR)拟合， $x_t$  均值一般为 0，设  $\hat{x}_t$  为  $x_t$  的估计值，则其自回归模型为

$$\hat{x}_t = \phi_0 + \phi_1 x_{t-1} + \phi_2 x_{t-2} + \dots + \phi_p x_{t-p} \quad (8)$$

其中， $p$  为模型阶数； $\phi_i$  为模型自回归系数， $i=0, 1, \dots, p$ 。

本文模型阶数的确定采用 AIC 准则，对参数的估计采用最小二乘法原理估计。

## 3 实例分析

为了直观地说明上述方法在动态关联规则的元规则挖掘中预测支持度值的具体过程，下面将针对一实例用基于时间序列模型方法对其支持度序列的值进行分析和预测。

实例：本文把文献[2]中的动态关联规则的算法做了修改后应用到某超级市场的2004年~2007年的销售数据库中。该数据库是SQL Server 2005自带的。挖掘出某规则  $A \Rightarrow B$ ，即表示，顾客在购买A的前提下购买B，它在2004年~2007年这4年内的每个月支持度计数(即顾客同时购买了A和B发生的次数，包含小于最小支持度计数的月份数据)为  $f_i$ ，对  $f_i$  建立时间序列模型，考察此规则的变化趋势，并对该规则在2008年~2010年这3年的支持度计数进行预测。

为了验证此模型的预测准确性，用此模型预测2004年~2006年的支持度计数的预测值，将2007年的预测值与实际值相比较后，可以说明所预测的2008年~2010年的数据是可靠的。

### 3.1 模型参数识别

依式(4)~式(8)进行计算，依次求取趋势项  $C_0 \sim C_7$  和  $R$ ，其结果分别为 242.962, -300.355, 149.358, -29.757, 2.993 6, -0.166 3, 0.005 2, -0.000 1 和 0.872 2。

周期项参数值如下：

傅氏系数	显著周期序号 $k$				
	0	3	4	6	8
$a_i$	-1.133 9e-010	-5.262 6	12.773 2	-3.938 6	-3.456 9
$b_i$		-3.554 8	9.678 9	-5.945 0	-6.129 4

随机项参数值如下：

$\phi_0$	$\phi_1$	$\phi_2$	$\phi_3$	$\phi_4$	$\phi_5$	$\phi_6$	$\phi_7$
$\phi_j$	-0.433 5	-0.868 3	-0.829 9	-0.823 6	-0.460 6	0.036 3	0.766 7
AIC	-769.725						

将趋势分量、周期分量和随机分量叠加即可得到支持度计数元规则预测模型，即

$$S_t = 242.962 - 300.355t + 149.358t^2 - 29.757t^3 + 2.9936t^4 - 0.1663t^5 + 0.0052t^6 - 0.0001t^7 + 0.8722t^8 - 3.5548\sin\frac{\pi t}{6} + 12.7732\cos\frac{2\pi t}{9} - 3.9386\sin\frac{2\pi t}{9} - 3.4569\cos\frac{\pi t}{3} - 5.9450\sin\frac{\pi t}{3} - 3.4569\cos\frac{4\pi t}{9} - 6.1294\sin\frac{4\pi t}{9} - 0.4335x_{t-1} - 0.8683x_{t-2} - 0.8299x_{t-3} - 0.8236x_{t-4} + 0.0363x_{t-5} + 0.7667x_{t-6} \quad (9)$$

### 3.2 精度检验

支持度计数元规则预测模型建立后，需对其精度进行检验。对实际值与预测值进行对比，拟合结果见图1。

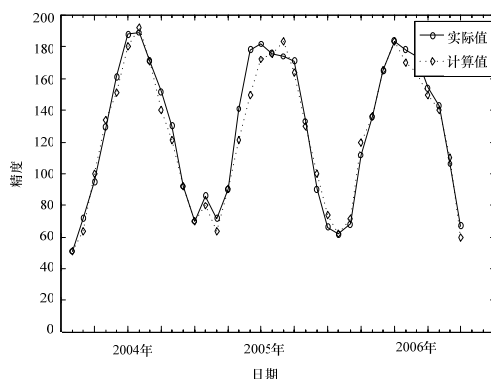


图1 预测值和实际值拟合曲线

在把已建立的模型用于预报前，也要进行精度检验，本

文采用2007年的数据进行后验预测检验，其预测结果见表1。

表1 2007年后验预测检验数据

时间	实际值	预测值	绝对误差	相对误差/(%)	时间	实际值	预测值	绝对误差	相对误差/(%)
1月	62	60	-2	3.20	7月	180	184	4	2.22
2月	82	79	-3	3.60	8月	178	174	-4	2.25
3月	107	105	-2	1.87	9月	154	160	6	3.90
4月	136	140	4	2.94	10月	142	140	-2	1.41
5月	165	170	5	3.03	11月	108	111	3	2.78
6月	182	180	-2	1.10	12月	66	69	3	4.55

经后验差计算检验，所建立的支持度计算元规则预测模型的后验差比值  $c$  和小误差频率  $p$  分别为：2.737 5 和 1。从图1中也可以看出其拟合精度较高。

### 3.3 预测分析

曲线拟合及后验预测检验表明此模型满足精度要求。依据式(9)，预测此规则2008年~2010年的每个月的支持度计数值，所得结果见表2。

表2 2008年~2010年逐月预测的支持度元规则数据

年份	1月	2月	3月	4月	5月	6月	7月	8月	9月	10月	11月	12月
2008年	52	60	105	120	164	186	175	179	142	128	105	70
2009年	56	84	98	133	180	188	187	175	160	143	107	61
2010年	53	119	126	163	160	195	177	171	161	143	97	50

从表2可见，利用本模型预测所得的此规则2008年~2010年逐月预测的支持度计数仍然呈波形按周期有规律地变动，同时每年的1月和12月的数据与往年的相比有所降低，可根据实际情况分析原因，进行决策分析，其他各月数据与往年基本保持一致。

## 4 结束语

本文利用基于时间序列的方法解决了动态关联规则过程中如何选取支持度的问题，克服了静态关联规则挖掘算法的数据集中各项差异较大或分布不均匀的缺点。并通过一个实例说明利用本模型挖掘动态关联规则的元规则的一般过程，结果表明，其计算简单，拟合精度和预测精度均较高，是一种较好的模拟预测模型。利用本模型进行元规则挖掘适用于历史数据较多且无规律可循时，若历史数据较少则可以使用灰色模型进行建模预测。通过对规则变化趋势的拟合和对未来数据的预测能更准确地把握规则的变化趋势，从而使动态关联规则挖掘在合理的元规则指导下得到更精确的结果。

## 参考文献

- [1] Liu Jinfeng, Rong Gang. Mining Dynamic Association Rules in Databases[C]//Proc. of International Conf. on Computational Intelligence and Security. Xi'an, China: [s. n.], 2005: 668-695.
- [2] 沈斌, 姚敏. 一种新的动态关联规则及其挖掘算法研究[J]. (2007-12-02). <http://www.paper.edu.cn>.
- [3] Han Jaiwei, Kamber M. 数据挖掘：概念与技术[M]. 范明, 孟晓峰, 译. 北京: 机械工业出版社, 2007.
- [4] 欧阳为民, 郑诚, 蔡庆生. 数据库中加权关联规则的发现[J]. 软件学报, 2001, 12(4): 612-619.
- [5] Wai Hoau, Chan K C C. Mining Changes in Association Rules: A Fuzzy Approach[J]. Fuzzy Sets and Systems, 2005, 14(1): 87-104.

编辑 张正兴