

18

时序数据库 知识发现 时序逻辑

计算机科学1999Vol. 26No. 8

68-70,89

# 基于时序逻辑的时序数据库中知识发现方法\*

Approaches of Discovering Knowledge in Temporal Databases Based on Temporal Logic

王清毅 范焱 蔡庆生 TP392 TP18

(中国科学技术大学计算机系 合肥230027)

**Abstract** This paper first gives the definitions of temporal database, temporal operator and temporal logic programming, and then discusses the temporal relations between the facts in the temporal databases. Some concepts such as fact sequence, class of temporal pattern are also defined. The approaches of discovering two classes of temporal patterns are described in detail.

**Keywords** Knowledge discovery in databases, Temporal databases, Temporal operator, Temporal logic programming, Temporal pattern

知识发现是一个众多学科诸如人工智能、机器学习、模式识别、统计学、数据库和知识库、数据可视化等相互交叉、融合所形成的一个新兴的且具有广阔应用前景的领域。目前国际上对知识发现的研究与开发进展很快<sup>[1]</sup>。众多现实世界(如天气预报、电信、金融)的数据库中数据都是随时间变化的,即数据具有时序性。目前,时序数据库中的知识发现问题正逐渐引起 KDD 研究者的兴趣。本文首先给出时序数据库、时序算子和时序逻辑程序的定义,然后讨论时序数据库中事实间的时序关系,定义了事实序列、时序模式的类等概念,详细描述时序数据库中两类时序模式的发现方法。为讨论方便,本文在同样的意义下不加区别地使用知识和模式这两个术语。

## 1 相关的概念

考虑和自然数集合  $N$  同态的时间域  $T$ , 一个关系由有限个非时序属性和一个或多个在时间域上变化的时序属性所构成。考虑数据库是由有限个谓词(即将一个元组视为一个谓词)构成的集合  $\zeta$ 。

**定义1** 一个时序数据库  $\zeta$  是一个函数, 该函数将  $\zeta$  中的每个谓词映射到  $T$  的一个子集(该子集是这个关系成立的时刻的集合), 即  $B = (\zeta \rightarrow 2^T)$ 。

在经典一阶逻辑中加进时序算子, 就构成了一阶时序逻辑(FOTL)。时序算子如下:

**定义2** 现在 always 算子  $\square$ 。  $\square A$  为真, 如果  $A$  在将来每个时刻为真。现在 sometimes 算子  $\diamond$ 。  $\diamond A$  为真, 如果  $A$  在将来某个时刻为真。next 算子  $\bigcirc$ 。  $\bigcirc A$  现在为真, 如果  $A$  在下一个时刻为真。过去 always 算子  $\blacksquare$ 。  $\blacksquare A$  为真, 如果  $A$  在过去每个时刻为真。过去 sometimes 算子  $\blacklozenge$ 。  $\blacklozenge A$  为真, 如果  $A$  在过去某个时刻为真。previous 算子  $\bullet$ 。  $\bullet A$  为真, 如果  $A$  在前一个时刻为真。until 算子  $\triangleright$ 。  $A \triangleright B$  为真, 如果  $A$  一直为真, 直到  $B$  出现, 此时  $A$  为假。该算子是一个二元算子。

约定  $\square, \diamond, \bigcirc$  的运算优先级高于  $\triangleright$  及  $\wedge, \vee, \rightarrow, \leftrightarrow$ , 而  $\triangleright$  的优先级又高于  $\wedge, \vee, \rightarrow, \leftrightarrow$ 。我们仍然沿用一阶逻辑中的项、原子公式、文字、子句以及基项、基原子公式、基文字和基子句的概念。一个时序文字就是一个以时序算子为前缀的文字。如 next-文字就是以有限个 next 算子为前缀的文字,  $\bigcirc \bigcirc \dots \bigcirc A = \bigcirc^n A$  表示在  $n$  个下一时刻  $A$  为真。一个时序规则具有 BODY  $\rightarrow$  HEAD 的形式, 其中 BODY 和 HEAD 都是时序文字的合取式, 合取式中各个时序文字用逗号分开。

一个时序逻辑程序 TLP 是一个时序规则的集合。这里我们采用  $Datalog_{it}$ <sup>[2]</sup> 来表示时序规则。在  $Datalog_{it}$  中, 每个谓词至多包含一个时序参数。另外我们允许谓词的否定出现在时序规则的头和体中。

\* )本文得到国家自然科学基金资助,王清毅 博士生,主要研究领域为机器学习、知识发现,范焱 博士生,主要研究领域为机器学习、知识发现,蔡庆生 教授、博士生导师,主要研究领域为人工智能、机器学习和知识发现。

允许时序变元之间的比较出现在时序规则的体中。将 $(\dots((0+1)+1\dots+1))$ (加  $n$  个 1) 记为  $n$ , 将 $(\dots((T+1)+1\dots+1))$ (加  $n$  个 1) 记为  $T+n$ , 例如, 存在这样一条规则: 对银行有欺诈行为的顾客将永远不允许在该银行存款, 用时序算子可以表示为:

$\diamond$ SAVING (bank, customer)  $\wedge$  SWINDLE (bank, customer)  $\rightarrow \square \neg$ SAVING (bank, customer)。

在 Datalog<sub>ts</sub> 中, 可以利用谓词中的时序参数来表示顾客的存款时间, 即用 SAVING (bank, customer, time) 来替换 SAVING (bank, customer), 用 SWINDLE (bank, customer, time) 来替换 SWINDLE (bank, customer)。

## 2 两个知识发现方法

一般说来, 时序数据库中的事实(即谓词的对象为常元)之间的时序关系可以形象地用有向非循环图如图 1 所示。

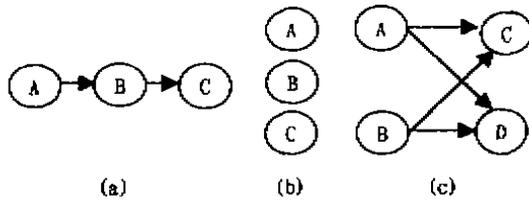


图 1 时序数据库中的事实间的时序关系

图 1(a) 描述了事实 A, B, C 之间的时序串行关系, 即 A 在 B 之前发生, B 在 C 之前发生; (b) 描述了事实 A, B, C 之间的时序并行关系, 即 A, B, C 的发生没有时间顺序上的关系约束; 在 (c) 中, 事实 C, D 只有在事实 A, B 发生以后才会发生, 但在 A 和 B 之间, C 和 D 之间则没有时间顺序上的约束关系。下面我们来形式化地定义时序数据库中的几个概念。

**定义 3** 给定一个基本事实类  $F$ , 一个时间域  $T$ , 一条事实  $fact$  是一个序偶  $(f, t)$ , 其中  $f \in F, t \in T$ 。一个事实序列  $S$  是一个三元矢  $(T, T', s)$ , 其中  $T$  是开始时间,  $T'$  是结束时间,  $s$  是一个事实的有序序列, 即  $s = (f_1, t_1), (f_2, t_2), (f_3, t_3), \dots, (f_n, t_n)$ , 其中  $f_i \in F, T_i \leq t_i \leq T', i = 1, 2, 3, \dots, n, t_i \leq t_{i+1}, i = 1, 2, 3, \dots, n-1$ 。

**定义 4** 一个时序模式的类  $P = (V, \leq, g)$ , 其中  $V$  是结点的集合,  $\leq$  是关于  $V$  的偏序,  $g$  是一个映射  $g: V \rightarrow F$ , 该映射将每个结点和一条事实相关

联,  $g(V)$  中的事实以  $\leq$  所确定的时间顺序发生。

**定义 5** 称一个时序模式的类  $(V, \leq, g)$  出现在一个事实序列  $s = (f_1, t_1), (f_2, t_2), (f_3, t_3), \dots, (f_n, t_n)$ , 如果存在一个满射  $h: V \rightarrow \{1, 2, \dots, n\}$ , 使得对所有的  $v \in V, g(v) = e_{h(v)}$ , 且对所有的  $v, w \in V$ , 有  $v \leq w, h(v) < h(w)$ , 称一个时序模式是频繁的, 如果这个时序模式在事实序列中发生的次数大于给定的一个频度阈值  $c$ 。

从时序逻辑公式的角度来看, 一个时序模式是一个基一阶时序逻辑公式。一个含有一个或多个变元的时序逻辑公式定义了一个时序模式的类。对时序逻辑公式中的变元用特定的基值例化定义了一个时序模式。

基于以上定义, 时序数据库中的知识发现任务可描述为: 给定一个事实序列, 一个时序模式的类及一个频度阈值, 从该事实序列中发现全部频繁的时序模式。下面我们讨论两类时序模式的发现方法。

### 2.1 一阶时序逻辑公式定义的时序模式类

我们先从如式 (1) 所示的时序模式类开始讨论发现方法, 然后再推广到一般情形。

$$a(X) \rightarrow \diamond b(Y) \wedge \bigcirc^{T+} b(Y) \triangleright c(Z) \quad (1)$$

其中  $a, b, c$  是时序数据库中的时序谓词,  $X, Y, Z$  分别是  $a, b, c$  的属性向量。该公式的直观意义是: 如果  $a$  现在为真, 那末最终会出现  $b$  为真的状态, 并且  $b$  将一直为真, 直到  $c$  出现的时刻, 此时  $b$  为假。问题是要发现全部  $X, Y, Z$  的例化即时序模式, 且这些时序模式是频繁的。图 2 是发现时序模式类 (1) 的 TLP 程序。我们在  $a(X), b(Y), c(Z)$  中分别引进了时序参数  $T$ , 还引进了一个时序识别谓词  $d(X, Y, Z, T)$ 。

- (1)  $\text{simtime}(0)$
- (2)  $\text{simtime}(T) \rightarrow \text{simtime}(T+1), \neg \text{simtime}(T)$
- (3)  $\text{simtime}(T), a(X, T) \rightarrow a(X, T-1) \rightarrow \text{flag1}(X, T1)$
- (4)  $\text{simtime}(T), a(X, T), a(X, T-1) \rightarrow \text{flag2}(X, T)$
- (5)  $\text{simtime}(T), b(Y, T), \neg b(Y, T-1) \rightarrow \text{flag3}(Y, T2), \neg \text{flag2}(X, T)$
- (6)  $\text{simtime}(T), b(Y, T), b(Y, T-1) \rightarrow \text{flag4}(Y, T)$
- (7)  $\text{simtime}(T), c(Z, T), \text{flag3}(Y, T2), \text{flag4}(Y, T-1), (T2 \leq T3 \leq T) \rightarrow d(X, Y, Z, T3)$
- (8)  $\text{simtime}(T), \neg b(Y, T), \text{flag3}(Y, T2), \text{flag4}(Y, T-1), (T2 \leq T3 \leq T) \rightarrow \neg \text{flag4}(Y, T3), \neg \text{flag3}(Y, T2), \neg \text{flag1}(X, T1)$
- (9)  $c(Z, T) \rightarrow d(X, Y, Z, T)$

图 2 发现式 (2) 所示时序模式的 TLP 程序

规则 (1) 以谓词  $\text{simtime}$  作为系统时钟来仿真时间; 规则 (2) 表示时间以“1”为时间粒度向前推进; 当谓词  $a(X)$  在第一个时间点为真时, 规则 (3) 即设置了  $\text{flag1}$  来标记; 规则 (4) 继续设置  $\text{flag2}$  来标记  $a$

(X)为真的时间段;当 b(Y)发生时,规则(5)设置了谓词 flag3,改变了 flag2;规则(6)继续设置 flag4来标记 b(Y)为真的时间段;当 c(Z)发生时,规则(7)设置了识别谓词 d(X,Y,Z),d(X,Y,Z)在从 flag3被设置的时刻起到 c(Z)发生的时刻止这一时间段为真;规则(8)在 b(Y)为真的时间段结束以后重新设置 flag4,flag3,flag1;规则(9)在 c(Z)为真时重新设置了 d(X,Y,Z)。

我们可以在以上的时序逻辑程序中加进一个计数器,对 q(X,Y,Z)为真的时间粒度进行计数,然后再将此计数器的值与所给的频度阈值进行比较,以确定出频繁的时序模式。将这一例子进行一般化推广便得到以下定理:

**定理** 在一个时序数据库 D 中,对任何一条一阶时序逻辑公式 φ 所确定的时序模式类,存在一个对应的时序逻辑程序 TLP 及一个识别谓词 d,使得对任何 D 中的有限事实,  $d \supseteq \varphi$ , 即 d 与 φ 同时成立。

该定理可以通过对 φ 中的时序算子的个数进行归纳来证明,其过程是构造性的。只要 D 有限且时序谓词的参数也有限, TLP 程序就会正常终止。

**2.2 由 next 算子  $\bigcirc$  和 and 算子  $\wedge$  定义的时序模式类**

这类时序模式如式(2)所示:  

$$P_0(X_0) \wedge \bigcirc^{k_1} P_1(X_1) \wedge \bigcirc^{k_2} P_2(X_2) \wedge \dots \wedge \bigcirc^{k_n} P_n(X_n) \quad (2)$$

其中,  $1 \leq k_1 \leq \dots \leq k_n$ ,  $P_0, P_1, \dots, P_n$  是不同的时序谓词,  $X_1, X_2, \dots, X_n$  是谓词的属性向量。

不难看出,该时序模式的频度是式(2)中自然数 n 的单调递减函数,即随 n 的增大,时序模式的频度逐渐降低,这样就一定存在一个 n 的最大值 N,此时的时序模式是非频繁的。我们的目标就是要找出这样的最大值 N,并且发现  $n < N$  时的所有频繁模式。实现这一目标的一个简单的方法就是从  $n=1$  开始,一步步随 n 以 1 递增来发现各个 n 值对应的频繁模式,每一步各自独立,不使用前次步骤的发现结果。换句话说,对第 i ( $1 \leq i < N$ ) 步的所有时序模式,不论频度高低,都要被计数。这样势必效率不高。考虑文[5]中讨论的发现数据库中关联规则的方法,在发现关联规则的时候,下一次迭代只利用前次迭代过程中生成的大的数据项集来生成新的候选数据项集以供进一步的评价。借鉴这一方法,在第 i 步发现频繁时序模式的时候,只利用第 i-1 步的发现结果,这就避免了对那些频度低的时序模式的计算,提高了算法的效率。为讨论方便,不妨先将式(2)简化成式

(3)所示的形式:  

$$P(X) \wedge \bigcirc^{k_1} P(X) \wedge \bigcirc^{k_2} P(X) \dots \wedge \bigcirc^{k_n} P(X) \quad (3)$$

发现该时序模式的 TLP 程序如图3所示:  
 $n=1, freq_0 = \text{所有频繁基谓词的 } P \text{ 的集合}$   
 repeat  
     通过时序逻辑程序 TLP<sub>n</sub> 计算  $freq_n$   
      $n=n+1$   
 until ( $freq_n = \emptyset$ ).  
 其中 TLP<sub>n</sub> 描述如下:  
 (1) simtime(0)  
 (2) simtime(T),  $T \leq T^* \rightarrow \text{simtime}(T+1)$   
 (3) simtime(T),  $P(X, T), P(X, T+K_1), \dots, P(X, T+K_{n-1}), freq_{n-1}(X, K_1, \dots, K_{n-1}), P(X, T+K_n), \text{count}(X, K_1, \dots, K_{n-1}, K_n, VAL) \rightarrow \text{holds}(T, X, K_1, \dots, K_{n-1}, K_n) \rightarrow \text{holds}(T, X, K_1, \dots, K_{n-1}, K_n), \text{count}(X, K_1, \dots, K_{n-1}, K_n, VAL+1)$   
 (4) simtime(T+1),  $\text{count}(X, K_1, \dots, K_{n-1}, K_n, VAL), (VAL > c) \rightarrow freq_n(X, K_1, \dots, K_{n-1}, K_n)$

图3 发现式(3)所示时序模式的 TLP 程序

$T^*$  是时间域的最大值,  $freq_n(X, K_1, \dots, K_{n-1}, K_n)$  是第 n 步所发现的全部频繁模式,  $\text{count}(X, K_1, \dots, K_{n-1}, K_n, VAL)$  是频度计算谓词, c 是给定的频度阈值,规则(3)中的谓词  $\text{holds}(T, X, K_1, \dots, K_{n-1}, K_n)$  避免了对频度的重复计算,其初始值为假。此程序以递归形式写出。它对每一个 i 值,只利用 i-1 步的发现结果,再发现 i 步中的频繁时序模式。程序一直执行到  $freq_n = \emptyset$  结束。

显然,此 TLP 程序可以推广到式(3)的一般情形,不同的地方只是程序中的谓词标记要复杂些,限于篇幅,不再详细给出。式(3)中的谓词组合是任意的,并且谓词的个数也是有限任意的,因此此程序存在着组合爆炸的危险。另外,程序的每次迭代都对应着一个 TLP 程序,所以程序的复杂度不再是线性的。

**结束语** 时序数据库中的知识发现是当前 KDD 领域一个较为热门的研究课题,有关论文尚不多见。本文介绍的方法是这方面一个初步尝试,以期起到抛砖引玉的作用。

**参考文献**

- 1 王清毅,陈恩红,蔡庆生. 知识发现的若干问题及应用研究 计算机科学,1997(5)
- 2 Baudient M, et al. Temporal Deductive Databases. Ch13 in Temporal Databases Theory Design and Implementation. Benjamin/Cummings, 1993
- 3 Abadi M, Manna Z. Temporal Logic Programming. J. of Symb. Computation, 1989, 8: 277~295
- 4 Tuzhulin A. Sim TL: A Simulation Language Based on

(下转第89页)