

# 基于粗集理论的知识自动获取方法

刘道华<sup>1,2</sup>, 原思聪<sup>1</sup>, 李湘英<sup>2</sup>, 王发展<sup>1</sup>

(1. 西安建筑科技大学机电学院, 西安 710055; 2. 信阳师范学院计算机科学系, 信阳 464000)

**摘要:** 分析了粗集理论的知识自动获取方法的基本原理和获取过程, 研究了数据记录规范化方法、属性归约算法、最小决策规则集的求解、规则提取方法, 并给出了自动获取方法的实例。实例证明了该算法的有效性。

**关键词:** 粗糙集理论; 知识自动获取; 专家系统开发工具; 属性归约; 数据规范化

## Method of Automatic Knowledge Gain Based on Rough Sets Theory

LIU Dao-hua<sup>1,2</sup>, YUAN Si-cong<sup>1</sup>, LI Xiang-ying<sup>2</sup>, WANG Fa-zhan<sup>1</sup>

(1. School of Mech. & Elec. Engineering, Xi'an University of Arch. & Technology, Xi'an 710055;

2. Dept. of Computer Science, Xinyang Normal University, Xinyang 464000)

**【Abstract】** This paper analyzes the basic principle of automatic knowledge gaining method based on rough sets theory, introduces the automatic gaining process, and studies the method of data generalization, the concrete algorithm of the attribute reduction, solution of the smallest policy-making rule set and the method of decision-making rule. A concrete example is given about the knowledge's automatic gaining method. The example shows the method is effective and feasible.

**【Key words】** rough sets theory; knowledge automatic gain; development tool of expert system; attribute reduction; data generalization

机械设计专家系统开发工具是当前机械设计制造业研究的热点, 作为开发具体专家系统的母板, 要提供多种知识自动获取方法, 以适合不同知识表示及推理的知识自动获取的需要, 知识的自动获取是人工智能知识自动获取研究的难点。在具体的应用规则库中, 通过一定的算法获得新的规则知识, 这是系统开发的难点。笔者通过数据规范化算法, 对规则库的记录集进行规范化, 通过具体知识规则库中的属性约简, 在对规则库辨识矩阵及辨识函数的求解中, 找出属性的核心, 从决策矩阵及决策函数中提取最大泛化的规则(最小决策规则集), 从属性归约及数据规范化后的规则库中, 提炼新知识规则。具体实例证明该方法在一定程度上解决了信息系统知识自动获取的“瓶颈”问题。

### 1 知识自动获取原理

基于粗糙集理论的基本集合运算法则, 通过对集合属性的辨识矩阵及辨识函数的计算, 找出规则前件的核属性, 通过对属性约简集中的非核属性的可靠度计算, 删除不必要的属性, 这样可大大简化规则的前件, 而对规则的个数通过信息表的记录元组的不相关性(不可重复性)进行约简, 通过对决策矩阵及决策函数的计算求解出最小决策规则集, 并通过概念提升技术对约简后的信息表进行规则提取。整个过程的每一个环节均可通过编程来自动实现。

### 2 粗集理论的知识自动获取过程

对知识的数据进行预处理, 并将其归一化为规范的数据库表, 对该信息表通过辨识矩阵及辨识函数求出属性的核心, 最终通过属性归约删除不重要的属性列, 对删除属性列后的信息表剔除重复的元组, 然后对元组再次进行数据规范化, 从

而提取更一般的元组, 最后通过计算决策矩阵及决策函数, 求解出最小决策规则集, 并从最小决策规则集中提取新知识规则存入知识库, 从而完成新的知识获取过程, 具体过程见图1。

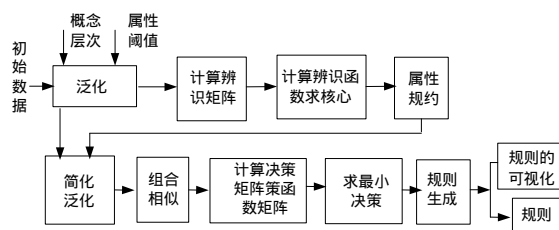


图1 粗糙集理论的知识自动获取过程

#### 2.1 数据规范化

实际数据库的元组数量往往很大, 但有许多元组的信息是冗余的; 有些元组数据虽然差异很少, 但在可忽视的范围内; 有的虽然不同, 但代表的是同一类信息, 可以等同视之。如果对这些数据不先进行处理或选择, 而直接用粗糙集算法进行属性归约, 效率势必很低。往往在使用归约算法前先将初始数据泛化, 缩减或浓缩数据库中的元组。数据泛化的常用方法是: 使用属性迁移和概念树提升技术, 逐个属性地进

**基金项目:** 国家自然科学基金资助项目(50275113); 陕西省教育厅产业化基金资助项目(04JC21)

**作者简介:** 刘道华(1974-), 男, 讲师、博士研究生, 主研方向: 智能系统, 优化设计; 原思聪, 教授、博士生导师; 李湘英, 硕士、副教授; 王发展, 博士后、副教授

**收稿日期:** 2006-11-21 E-mail: ldh30708@sina.com

行泛化,最后得到一个泛化关系<sup>[1]</sup>。

在基于粗糙集理论的数据归约中,元组归约和属性归约都是针对信息系统进行的,数据泛化的目的是生成泛化的信息系统,以便下一步针对该信息系统进行属性归约。

下列算法用于从一个初始信息系统生成一个泛化的信息系统。其中,  $d_i$  表示关系  $R$  中某个属性  $a_i(i=1,2,\dots,n)$  的不同值的个数;  $t_i$  表示“属性阈值”;  $p$  表示支持泛化的全部元组在整个数据库中所占的百分比(支持率),泛化时由用户指定。如果一个属性是可泛化的(即存在相应的概念层次树  $H_i$ ),并且  $d_i > t_i(t_i > 1)$ ,则称该属性相对于  $t_i$  是可泛化的。该算法的特点在于:用  $d_i/t_i$  值选择下一个被泛化属性,保证它具有最大  $d_i/t_i$  值,这样可以最大程度地改善泛化性能。

**算法** 从某个关系中抽取泛化的信息系统。

**输入**

- (1)与任务相关的数据集  $R$ , 属性  $a_i(i=1,2,\dots,n)$ ;
- (2)概念树集合  $H$ ,  $H_i \in H$  是可泛化属性  $a_i$  的概念层次;
- (3)属性  $a_i$  对应的属性阈值  $t_i$ , 以及  $p(0 < p \leq 1)$ 。

**输出** 泛化的信息系统  $R'$ 。

**步骤**

- (1)  $\max tuples \leftarrow p \times |R|, R' \leftarrow R$ ;
- (2)计算每个属性的  $d_i$  值;
- (3)While  $|R'| \geq \max tuples$  and  $\exists d_i > t_i$

Do {选择属性  $a_i$ , 使得  $a_i \in A$ , 并满足  $d_i/t_i$  最大;

If 属性  $a_i$  是可泛化的,

Then 在  $H_i$  中将  $a_i$  的概念提升一级,并在  $R'$  中对  $a_i$  下的概念进行替换,

Else 从  $R'$  中删除属性  $a_i$ ,

EndIf

删除  $R'$  中的重复元组,重新计算每个属性的  $d_i$  值},

End while

利用上述方法对数据集进行处理,可以大大提高数据归约的效率。

## 2.2 属性归约

基于粗糙集理论的属性归约的一般步骤为:

- (1)通过辨识矩阵及辨识函数求出属性归约集的核心;
- (2)运用归约算法计算归约集,并根据归约集属性依赖度的计算删去不重要的属性,当数据量很大时,可在求核心前先将初始数据表泛化。

### 2.2.1 通过辨识矩阵和辨识函数求解属性的核心

属性的核心可作为计算归约集的基础,属性的核心是由辨识矩阵内仅含有单个元素构成的集合组成<sup>[2]</sup>。

令  $S=(U,A,V,f)$  是一个知识表达系统,  $|U|=n$ ,  $S$  的辨识矩阵  $M$  是一个  $n \times n$  矩阵,其任一个元素为

$$\alpha(x_i, x_j) = \begin{cases} \Phi & x_i, x_j \in D \text{ 的同一等价类} \\ \{a \in A | f(x_i, a) \neq f(x_j, a)\} & x_i, x_j \in D \text{ 的不同等价类} \end{cases} \quad (1)$$

因此,  $\alpha(x_i, x_j)$  是区别对象  $x_i$  和  $x_j$  的所有属性的集合。

为了求得辨识矩阵中的核,引入一个布尔函数,称其为辨识函数(discernible function),用  $\Delta$  表示,对每个属性  $a \in A$ ,指定一个布尔变量“ $a$ ”。

若  $\alpha(x_i, x_j) = \{a_1, a_2, \dots, a_k\} \neq \Phi$ , 则指定一个布尔函数  $a_1 \vee a_2 \vee \dots \vee a_k$ , 用  $\sum \alpha(x_i, x_j)$  来表示;若  $\alpha(x_i, x_j) = \Phi$ , 则指定布尔常量为 1。辨识函数  $\Delta$  可定义为

$$\Delta = \prod_{(x_i, x_j) \in U \times U} \sum \alpha(x_i, x_j) \quad (2)$$

由式(2)可知,利用布尔代数的分配律、结合律及吸收律求得该函数的“极小析取范式”的所有“合取式”,其“极小范式”为属性集  $A$  的所有约简。

约简是满足能区别由整个属性集区别的所有对象的属性极小子集。如果  $B \subseteq A$  满足条件  $B \cap \alpha(x_i, x_j) \neq \Phi$ ,  $\forall \alpha(x_i, x_j) \neq \Phi$  的极小子集,则  $B$  是  $A$  的一个约简。

而核是辨识矩阵中所有单个元素组成的集合,即

$$core(A) = \{a \in A | \alpha(x_i, x_j) = \{a\}\} \quad (3)$$

其中,  $x_i, x_j \in U$ 。

根据定理,可以通过辨识矩阵计算出核心<sup>[3]</sup>。

**定理** 设有信息系统,即

$$S = \{U, Q, V, f\}, \quad Q = C \cup D$$

对于任何  $c \in C, c \in core(C, D), M(C) = (m_j)$ , 当且仅当存在  $1 \leq j < i \leq n, m_j = \{c\}$ 。即  $core = \{c \in C | m_j = \{c\}, 1 \leq j < i \leq n\}$ 。

### 2.2.2 最简属性归约集的求解

在决策过程中,每个归约集都可以代替整个条件属性集,而不改变原有的依赖关系。通常情况下,一个信息系统可能有多个属性归约集<sup>[4]</sup>。但每一个归约集必含有属性核,但除核属性之外的属性如何删除,可采用的方法为:先求得等价类的依赖度,然后计算除去约简核属性之外的属性的依赖度,从而作出是否删除的判断。

在属性归约中,利用 2 个属性集合  $P, R \subseteq Q$  之间的相互依赖程度,可以确定某一属性  $a$  的重要性。属性集  $P$  对  $R$  的依赖程度用  $\gamma_R(P)$  表示。定义为

$$\gamma_R(P) = \frac{card(POS_R(P))}{card(U)} \quad (4)$$

$$POS_R(P) = \bigcup_{X \in U/IND(P)} apr_R(X) \quad (5)$$

式(4)中,  $card(\cdot)$  表示集合的基数;式(5)中,  $POS_R(P)$  是属性集  $R$  在  $U/IND(P)$  中的正区域。而对于任何一个属性集合  $P \subseteq Q$ , 不可分辨关系  $IND$  为

$$IND(P) = \{(x_i, x_j) \in U \times U | f(x_i, a) = f(x_j, a), \forall a \in P\} \quad (6)$$

而属性集  $C \subseteq A$  对论域  $U$  进行划分得到的等价关系族记为  $U/C$ 。对任何一个对象子集  $X \subseteq U$  和属性子集  $R \subseteq Q$ ,  $R$  的“下近似”定义为

$$apr_R(X) = \bigcup \{Y \in U/IND(R) | Y \subseteq X\} \quad (7)$$

通过属性依赖度的求解,计算约简集中某一属性的重要性,来决定是否删除某一属性,考察条件属性  $a$  对于决策  $d$  的重要性,该指标可以用删除条件属性  $a$  前后的重要性之差  $\beta$  来表示。

$$\beta = \gamma_R(P) - \gamma_{R-\{a\}}(P) \quad (8)$$

其中,若  $\beta = 0$ , 则表明条件属性  $a$  对于决策  $d$  没有影响,可以在系统中将其删除,经过分析处理之后,能够得到决策系统的约简。

属性归约后并不改变原始属性集与决策属性之间的依赖程度,如果去掉该属性会造成依赖度变化,则恢复该属性,否则剔除该属性<sup>[5]</sup>。最后剩下的属性集就是最佳归约集或是用户定义的最小属性集。

### 2.3 最大泛化规则(最小决策规则集)的求解

对于决策属性  $d \in A$ , 及其特定值  $V_d$ , 关注的是满足  $d(e) = V_d$  的对象  $e$  的集合  $\{V_d\}$ 。用矩阵形式可以将区分所有属于集合  $\{V_d\}$  的对象和属于  $U - \{V_d\}$  集合的对象的属性值对表示出来<sup>[6]</sup>。

