

# 基于粗糙集的确定性控制规则决策模型

邓九英<sup>1,2</sup>, 毛宗源<sup>1</sup>, 杜启亮<sup>1</sup>, 谭光兴<sup>1</sup>

(1. 华南理工大学自动化科学与工程学院, 广州 510641; 2. 广东教育学院计算机科学系, 广州 510303)

**摘要:**粗糙集属性分区数的变化会影响属性重要性和属性对决策属性的支持度。该文对知识表示系统的数据相关性进行分析, 综合考虑系统的泛化能力, 提出能生成确定性控制规则的决策模型, 给出决策模型中属性分区数求取以及属性相对约减产生的判断与算法实现。实验结果表明, 该算法简洁有效, 验证了决策模型的准确性与实用性。

**关键词:**决策模型; 确定性控制规则; 属性分区; 属性重要性; 支持度

## Decision Model of Certain Control Rules Based on Rough Set

DENG Jiu-ying<sup>1,2</sup>, MAO Zong-yuan<sup>1</sup>, DU Qi-liang<sup>1</sup>, TAN Guang-xing<sup>1</sup>

(1. College of Automatic Science & Engineering, South China University of Technology, Guangzhou 510641;

2. Department of Computer Science, Guangdong Institute of Education, Guangzhou 510303)

**【Abstract】** Different attributes partitions affect about attribute significance and attributes support degree to decision attribute. This paper analyzes the relationship among data in knowledge represent system based on rough set. Generalization of system is combined to consider. A decision model is presented for generating certain control rules. Judgment and algorithm are introduced to compute attributes partitions in decision model and find relative attributes reduction in decision model. Experimental results show conciseness and efficiency of algorithms and preciseness and utility of decision model.

**【Key words】** decision model; certain control rules; attribute partition; attribute significance; support degree

### 1 概述

粗糙集方法是人工智能(AI)和认知学的基础, 也是机器学习、知识获取、决策分析和从数据库中发现知识等领域的基本工具, 在决策支持系统和数据挖掘方面粗糙集显得格外重要。在粗糙集方法用于数据分析时, 数据模式特征用近似的方法表示, 或等价于从数据中归纳出的决策规则<sup>[1]</sup>。科研人员围绕决策规则做了大量的研究工作, 包括用粗糙集评价重要规则的方法<sup>[2]</sup>、选择具有高特征的属性约减子集算法和获取具有泛化能力的简洁规则策略<sup>[3]</sup>等。

通过知识获取生成决策规则的质量, 与规则应用到实践中的控制结果的好坏有直接关系<sup>[4]</sup>, 对生成与选择决策规则质量的评判依据包括几方面, 其中主要的一点是规则的准确性; 另外规则的简洁性直接关系到规则的实用性。由于属性的分区数(量化值的个数)变化会影响属性的重要性和属性对决策属性的支持度, 因此本文提出一种构建确定性决策规则的模型。

### 2 决策表

#### 2.1 定义与性质

决策表可用来定义知识表示系统(Knowledge Representation System, KRS), 设 $S=(U, A)$ 为知识表示系统,  $C, D \subset A$ 是2个属性子集, 分别为条件属性与决策属性; KRS也可表示为 $S=(U, A, C, D)$ 。关系 $IND(C)$ 和 $IND(D)$ 的等价类分别称为条件等价类与决策等价类<sup>[5]</sup>。其中, 对于 $B \in A, IND(B) = \{(x, y) \in U \times U: \text{对任意} a \in B, a(x) = a(y)\}$ 。

对任意 $x \in U$ , 存在函数 $A \rightarrow V$ , 以及对任意 $a \in C \cup D$ , 有 $des_x(a) = a(x)$ ; 则函数 $des_x$ 称作( $S$ 中的)决策规则。限于对 $C$ 的

$des_x$ (表示为 $des_x|C$ )和限于对 $D$ 的 $des_x$ (表示为 $des_x|D$ )分别称为 $des_x$ 的条件与决策(或动作)。

如果对任意 $y \neq x, des_x|C = des_y|C$ 蕴涵 $des_x|D = des_y|D$ , 那么决策规则 $des_x$ 是在( $S$ 中)一致的; 否则, 决策规则是不一致的。如果所有的决策规则是一致的, 那么决策表是一致的; 否则, 决策表是不一致的。如果决策表是不一致的, 只要删除 $x$ (或 $des_x$ )就可把决策表变为一致。

#### 2.2 决策属性的支持度

为了计算某个属性(或属性集)的重要性, 常见的方法是: 从表中去除这个属性, 看属性去除前后分类的变化情况, 如果属性删除后的分类变化大, 就说明这个属性的重要性高; 反之, 这个属性的重要性低。

**定义 1** 假设一给定知识库 $S=(U, R)$ , 子集 $X \subseteq U$ 和等价关系 $R \text{ IND}(K)$ (缩写为 $R \text{ K}$ ),  $R$ -下近似集表示为:  
 $RX = \{Y \in U/R, Y \subseteq X\}$  or  $RX = \{x \in U, [x]_R \subseteq X\}$ 。

另一种表示形式为:  $POS_R(X) = RX, X$ 的 $R$ -正域。

**定义 2** 决策属性 $y \in D$ , 决策子集 $W \subseteq U/y$ , 决策属性 $y$ 关于条件属性 $a \in C$ 的支持子集定义为

$$S_a(y) = \bigcap_{u \in y} POS(U/a) = \bigcap_{u \in y} \bigcap_{v \in u/a, v \in W} v \quad (1)$$

$y$ 关于 $a$ 的支持度表示为 $spt_a(y) = |S_a(y)| / |U|$ 。

令 $W \subseteq U$ 为 $U$ 的子集, 决策表 $(U, C \cup D)$ ,  $W$ 关于条件属性

**作者简介:**邓九英(1962-), 女, 副教授, 主研方向: 人工智能, 数据挖掘技术, 仿真技术; 毛宗源, 教授、博士生导师; 杜启亮、谭光兴, 博士研究生

**收稿日期:** 2007-06-28 **E-mail:** djy1111@126.com

集  $X \subseteq C$  的支持子集可表示为  $S_X(W) = \bigvee_{U/X, V \subseteq W} V$ 。  
 $spt_X(W) = |\bigwedge_{U/X} POS(U/X) \cap W| / |U|$  称为  $W$  关于  $X$  的支持度。

### 3 核属性与决策规则

设决策表  $S=(U, A, V, f)$ ,  $A=C \cup D$ ,  $C \cap D = \emptyset$ 。令  $X_i$  和  $Y_j$  分别代表  $U/C$  与  $U/D$  中的等价类,  $des_{X_i}$  是等价类  $X_i$  的描述, 即等价类  $X_i$  对于各条件属性值的特定取值;  $des_{Y_j}$  是等价类  $Y_j$  的描述, 即等价类  $Y_j$  对于各决策属性值的特定取值<sup>[6]</sup>。

**定义 3** 决策规则定义为  $r_{ij} : des_{X_i} \rightarrow des_{Y_j}$ ,  $Y_j \cap X_i \neq \emptyset$ 。规则的确定性因子:

$$\mu(X_i, Y_j) = |Y_j \cap X_i| / |X_i|, 0 < \mu(X_i, Y_j) \leq 1 \quad (2)$$

如果  $\mu(X_i, Y_j) = 1$ , 则  $r_{ij}$  是确定的; 如果  $0 < \mu(X_i, Y_j) < 1$ , 则  $r_{ij}$  是不确定的。在得出决策规则之前, 决策表需要进行属性约减。

假设知识表示系统  $S=(U, A, V, f)$ ,  $n=|U|$ 。系统  $S$  的区分矩阵是一个  $n \times n$  维的矩阵, 其中的元素为  $a(x, y) = \{a \in A, |f(x, a) \neq f(y, a)\}$ 。如果  $a(x, y) = \{a_1, a_2, \dots, a_k\} \neq \emptyset$ , 用符号  $\Sigma a(x, y)$  表示布尔函数  $a_1 \vee a_2 \vee \dots \vee a_k$ ; 如果  $a(x, y) = \emptyset$ , 则取值为布尔常量 1。区分函数  $\Delta$  定义为

$$\Delta = \bigwedge_{(x, y) \in U \times U} \Sigma a(x, y) \quad (3)$$

如果  $B \subseteq A$  是一个满足条件  $B \cap a(x, y) \neq \emptyset, \forall a(x, y) \neq \emptyset$  的最小小子集, 则  $B$  就是  $A$  的一个约减。

核属性是由区分函数中单个因子组成的集合, 表示如下:

$$core(A) = \{a \in A, |a(x, y) = \{a\}, x, y \in U\} \quad (4)$$

### 4 确定性决策规则模型

为了使得到的决策规则简短而又准确, 构建 KR-system 决策模型的主要任务包括 2 个方面: (1) 确定属性的分区数; (2) 进行属性约减, 在满足相对约减子集对决策属性的支持度趋于 1 (或等于 1) 的条件下, 使相对约减子集中的元素个数较少。

由于属性的量化分区数的不同会改变属性的重要性参数, 如果要生成确定性决策规则, 应确定使属性对决策属性的支持度为用户给定数值的临界(最小)分区数<sup>[4]</sup>, 属性分区数的取值应大于这个临界值。从系统的泛化能力考虑, 属性约减子集中的元素个数应大于某基数(视系统的具体情况而定)。本文在此基础上提出 2 个判据和相应的 2 个算法, 分别用于确定决策模型中的属性分区数和属性约减子集。

假定知识表示系统不包含冗余与不相容数据, 属性取值为连续量, 对具有固定明确状态值的属性, 不适合用以下的判据与算法。

**判据 1** 属性分区的确定

- (1) 选取每个条件属性  $a_i \in C$  的分区数临界值,  $divnum(a_i)$ ;
- (2) 设定条件属性  $C$  的分区数初值,  $divC = \max\{a_i, |a_i \in C\}$ ;
- (3) 设定决策属性  $D$  的分区数初值,  $divD$ ;

(4) 循环执行,  $divC+1 \rightarrow divC, divD+1 \rightarrow divD$ , 属性分区量化取值后, 求  $Core$  中的元素个数  $sum$ , 直至  $sum$  大于给定基数;

(5) 循环执行,  $divC+1 \rightarrow divC, divD+1 \rightarrow divD$ , 属性分区量化取值后, 计算区分函数, 直至区分函数中只包含核属性因子;

(6) 后推计算,  $divC-1 \rightarrow divC, divD-1 \rightarrow divD$ , 属性分区量化取值后, 计算区分函数, 直至区分函数中出现了除核属性外的其他属性因子;

(7) 停止。

得出的  $divC$  和  $divD$  分别是系统  $S$  中条件属性与决策属性的分区数。

**判据 2** 属性相对约减子集的确定

(1) 设定属性约减子集  $X$  为属性核, 计算  $X$  对决策属性的支持度  $spt(X)$ ;

(2) 在几个约减方案中选择  $a_j$ , 并计算  $a_j+X$  对决策属性的支持度  $spt(a_j+X)$ , 使  $spt(a_j+X)-spt(X)$  为最大的  $a_j$  加入  $X$ ;

(3) 重复(2), 直至  $spt(X)=1$ ;

(4) 停止。

在 KRS 具有约减属性及其明确的属性数字量后, 应用粗糙集理论的决策算法, 可从 KRS 中抽取决策规则, 再对决策规则进行约减, 得出实用的精简决策规则。

## 5 决策模型的算法实现

### 5.1 决策模型的相应算法

假定在系统  $S=(U, A, V, f)$  的数据预处理阶段, 完成了每个属性  $a_i \in C$  对决策属性有极大支持度的最小属性分区数  $divnum(a_i)$  的计算, 运用算法 1 可计算出确定性规则的最小属性分区数, 运用算法 2 可计算出系统的属性约减的合理子集。

引入一些符号表示参数:  $n=Card(C)$ ,  $m=Card(Core)$ ,  $Ns=|U|$ ,  $sum=3$  (泛化能力基数),  $Card(D)=1$ 。

**算法 1** 计算属性分区数(AAP)

AAP 1: 赋初值,  $divnum(i), i=1, 2, \dots, n$ ;

AAP 2:  $divC = \max\{divnum(i), i=1, 2, \dots, n\}$ ;

AAP 3: 赋初值,  $divD$  (比如 4),  $sum$  (比如 3),  $m=1$ ;

AAP 4: Do while  $m \leq sum$

{ $divC+1 \rightarrow divC, divD+1 \rightarrow divD$ , 计算  $Core, Card(Core) \rightarrow m$ };

AAP 5: Do while  $m < Card(\Delta)$

{ $divC+1 \rightarrow divC, divD+1 \rightarrow divD$ , 计算  $\Delta$  与  $Core, Card(Core) \rightarrow m$ };

AAP 6: Do while  $m = Card(\Delta)$

{ $divC-1 \rightarrow divC, divD-1 \rightarrow divD$ , 计算  $\Delta$  与  $Core, Card(Core) \rightarrow m$ };

AAP 7: End

将系统  $S$  按照分区数  $divC$  和  $divD$  对属性进行量化取值。

**算法 2** 求属性相对约减子集(AADS)

AADS 1:  $Core \rightarrow X$ , 计算  $spt(X)$ ;

AADS 2:  $K = Card(\Delta) - Card(X)$ ;

AADS 3: Do while  $spt(X) < 1$  且  $K > 0$

{对所有  $\{\Delta-X\} a_j, j=1, 2, \dots, K, X+a_j \rightarrow X'$ , 找出  $spt(X') - spt(X)$  值最大的  $a_j$ ;

$X+a_j \rightarrow X$ ; 计算  $spt(X)$ };

AADS 4: End

这时, 只包含  $X$  属性子集的系统  $S$ , 可用于归纳出决策规则。

### 5.2 实验结果

选取广州某化工厂锌铝白煅烧过程实测数据, 采样时间为 2004-2-3 12:35 PM 至 2004-2-7 10:33 AM, 每 10 min 采样一次数据, 采样数据一共分为 19 个属性。剔除异常与无关数据, 保留了 13 个属性与 200 组数据, 其中, 条件属性用符号表示:  $C_1$  为窑头温度(ytwd),  $C_2$  为干燥温度(gzwd),  $C_3$  为排风温度(pfwd),  $C_4$  为#1 进料电流(jld1),  $C_5$  为#1 进料频率反馈(jlfreq1),  $C_6$  为#2 进料电流(jld2),  $C_7$  为#2 进料频率反馈(jlfreq2),  $C_8$  为#3 进料电流(jld3),  $C_9$  为#3 进料频率反馈(jlfreq3),  $C_{10}$  为#4 进料电流(jld4),  $C_{11}$  为#4 进料频率反馈(jlfreq4),  $C_{12}$  为立德粉粉种(fenkind); 决策属性,  $d$  为煅烧温度(ds wd)。

决策属性  $d$  的分区数初始值为 5, 经过数据预处理, 得出条件属性对决策属性支持度给定数值时的临界分区数如表 1 所示。

(下转第 170 页)