

文章编号: 1000-8152(2007)02-0183-04

一种基于有序属性决策系统分类规则提取策略

张文字

(西安邮电学院 管理系, 陕西 西安 710061)

摘要: 分类规则的精度取决于分类算法的构造。论文在综合分析基本粗糙集合概念及其约简算法的基础上, 阐述了一种基于准则的有序属性决策系统的数据挖掘算法。为此首先介绍了基于有序属性决策系统的集合表达, 然后利用有序属性决策系统中准则集与属性集的基本特征构造上下近似扩展模型, 得到准则集决策系统的四个相关参数。并进一步提出相应的数据约简与分类规则提取算法。最后给出了用此算法约简有序属性决策系统的算例, 实验结果表明此方法挖掘出的规则简练, 更具合理性和可靠性。

关键词: 数据挖掘; 粗糙集; 决策系统; 准则; 分类质量

中图分类号: TP181 文献标识码: A

Classification rule extracting strategy based on decision system with ordered attributes

ZHANG Wen-yu

(Department of Management, Xi'an Post and Telecommunication College, Xi'an Shaanxi 710061, China)

Abstract: The precision of classification rule is decided by the construction of classification algorithm. By the concepts and attribute reduction algorithm of basic rough set, a data mining algorithm based on the ordered character of attribute in decision system is proposed in this paper. First, the aggregation expression in decision system with ordered character of attribute is briefly introduced. Then, based on the basic characterization of criteria sets and attribute sets in decision system with ordered attributes, the upper and lower approximation expansion models are constructed to obtain the four relative parameters in decision system with ordered attributes. Thirdly, the corresponding data mining and classification rule extracting algorithm is constructed by using the proposed approach. Finally the rationality of the ordered attribute reduction method is validated by simulation example, and the result shows the rules mined by the method are concise and reliable.

Key words: data mining; rough set; decision system; criteria; classification quality

1 引言(Introduction)

数据挖掘是当前人工智能和数据库技术的研究热点之一。目前, 有许多基于机器学习、模式识别及统计学的规则获取方法, 近年来的粗集理论也为数据挖掘提供了重要工具, 其导出的规则精炼而且便于存储和使用。粗糙集是波兰数学家Z.Pawlak在1982^[1,2]年提出的一种分析数据的数学理论, 其特点是不需要预先给定某些特征或属性的数量描述, 而是直接从给定问题的描述集合出发, 依据不同的观测点把对象集划分成等价类来确定给定问题的近似域, 从而找出该问题中的本质特征和内在规律^[3]。

本文在简述基本粗集基本概念的基础上, 采用属性的有序特征对决策系统进行简化进而提取规

则, 并利用算例对该方法进行了合理性验证。

2 基本粗糙集(Basic rough set)

定义1 称 $S = (U, Q, \{V_q\}, f)$ 为知识表示系统。其中: U 为非空有限集, 称论域; Q 为非空有限集, 称属性集合; V_q 为属性 $q \in Q$ 的值域; $f : U \rightarrow V_q$ 为一单射, 使论域 U 中任一元素取属性 q 在 V_q 中的某一唯一值。如果 Q 由条件属性集合 C 和决策属性集合 D 组成, C, D 满足 $C \cap D = \emptyset, C \cup D = Q$, 则称 S 为决策系统。

定义2 令 $S = (U, C \cup \{D\})$ 为一决策系统, $B \subseteq C$, 定义 B 的不可分辨关系 $IND(B)$ 为 $IND(B) = \{(x, x') \in U^2 | \forall a \in B, a(x) = a(x')\}$, $a(x)$ 为元素 x 在属性 a 上的值。若 $(x, x') \in IND(B)$,

说明根据已有的信息不能将 x 和 x' 区分开。

定义3 令集合 V_d 表示决策属性 d 的域值,不失一般性,设 $V_d = \{1, 2, \dots, r(d)\}$,则决策属性 d 将论域 U 划分成 $\{Y_1, Y_2, \dots, Y_{r(d)}\}$,其中 $Y_k = \{x \in U | d(x) = k, 1 \leq k \leq r(d)\}$, $d(x)$ 为决策属性 d 在对象 x 上的值,集合 Y_k 称为决策系统 S 的第 k 个决策类或决策概念。

利用基本粗集理论对决策系统进行数据约简目前已经有很多的方法并在实际应用中取得了可喜的成果。例如基于正域的约简方法、基于分辨矩阵的约简方法、基于信息熵的约简方法、基于距离函数的约简方法^[4~11]等,每种方法都具有各自的特点及应用场合。下面在Salvatore Greco, Roman Slowinski^[12]等教授提出的基于准则决策系统的数据约简的方法基础上提出带有有序属性决策系统的数据挖掘新算法,并将其应用于电信网络质量评价中。

3 带有有序属性的粗糙近似(Rough approximation with ordered attributions)

3.1 基本原理(Basic theory)

基于有序属性的数据挖掘方法同样应用于如定义1所示的决策系统中,但是决策系统中的属性是通过某种递增递减顺序进行排序而形成的准则域。对于 $\forall q \in C$,且 q 为一准则,则存在优先关系 S_q ,使 xS_qy 表示“针对属性 q 而言, x 至少与 y 一样好”,则 S_q 是定义在 U 上的完备传递二元有序关系。相反,对于 $\forall q \in C$, q 不是准则而是属性,则存在等价关系 I_q ,使 xI_qy 为基本粗集理论中的等价二元关系,即它满足自反性、对称性和传递性。标记 $C^>$ 为 C 中的准则子集,而 $C^=$ 为 C 中的属性子集即 $C^> \cup C^= = C$, $C^> \cap C^= = \emptyset$ 。而且对于 $\forall P \subseteq C$,标记 $P^>$ 为包含在 C 中的准则集,即 $P^> = P \cap C^>$, $P^=$ 为包含在 C 中的属性集,即 $P^= = P \cap C^=$,且满足 $P = P^> \cup P^=$, $P^> \cap P^= = \emptyset$ 。令 R_p 是一个 U 上的自反和传递二元关系,即 R_p 是在 U 上部分有序的关系,精确地说,对于 $\forall P \subseteq C$,可定义 R_p 如下: $\forall x, y \in U, xR_p y$,当且仅当若 $q \in P^>$ 则存在 xS_qy 关系,若 $q \in P^=$,则存在 xI_qy 。如果 $P \subseteq C^>$ 及 $xR_p y$,则表示对于 $\forall q \in P$, x 优先于 y 。广义地讲,令 $CL = \{CL_t, t \in T\}$, $T = \{1, 2, \dots, n\}$,是 U 上的一个划分集,即对于 $\forall x \in U$, x 属于且仅属于一个 $CL_t \in CL$ 。假定 $\forall r, s \in T$,若 $r > s$,则 CL_r 的元素优于 CL_s 的元素。 CL 类的划分代表 U 上对象的综合评价:最差的对象位于 CL_1 中,最好的对象位于 CL_n 中,其他的对象属于 CL_r 。

从 CL 划分出发,可定义如下集合:

$$\begin{cases} CL_t^{\geq} = \bigcup_{s \geq t} cl_s, \\ CL_t^{\leq} = \bigcup_{s \leq t} cl_s. \end{cases}$$

对于 $\forall P \subseteq C$,令

$$\begin{cases} R_p^+(x) = \{y \in U : yR_p x\}, \\ R_p^-(x) = \{y \in U : xR_p y\}. \end{cases}$$

对于 $x \in U$, $R_p^+(x)$ 表示在准则集中优先于 x 的对象与在属性集 P 中与 x 等价的对象之和。类似地 $R_p^-(x)$ 表示在准则集 P 中被 x 优先的对象与在属性集 P 中与 x 等价的对象之和。

3.2 基于准则的上下近似扩展模型(Upper and lower approximation expanding model based criteria)

假设 P_p 为自反的,可以说 $x \in U$ 确切属于 CL_t^{\geq} 等价于 $R_p^+(x) \subseteq CL_t^{\geq}$ 。对于 $P \subseteq C, t \in T$, $y \in U$ 可能属于 CL_t^{\geq} 等价于至少存在一个对象 $x \in CL_t^{\geq}$ 使得在准则集 P 中 y 优先于 x 及在属性集 P 中 y 等价于 x ,即 $y \in R_p^+$ 。因此考虑 $P \subseteq C$,属于 CL_t^{\geq} 的所有对象集,确切地构成了 CL_t^{\geq} 的下近似,而可能属于 CL_t^{\geq} 的所有对象集构成了 CL_t^{\geq} 的上近似。

形式描述如下: $\forall t \in T$ 及 $\forall P \subseteq C$,可定义针对 P 而言 CL_t^{\geq} 的下近似 \underline{PCL}_t^{\geq} ,及针对 P 而言 CL_t^{\geq} 的上近似 \overline{PCL}_t^{\geq} :

$$\begin{cases} \underline{PCL}_t^{\geq} = \{x \in U : R_p^+ \subseteq CL_t^{\geq}\}, \\ \overline{PCL}_t^{\geq} = \bigcup_{x \in CL_t^{\geq}} R_p^+(x). \end{cases}$$

同理可得下近似 \underline{PCL}_t^{\leq} 与上近似 \overline{PCL}_t^{\leq} 的数学表示。

类似基本粗集定义,本文提出了基于准则集决策系统的相关参数,并将其应用在带有准则决策系统的数据约简算法中。

参数1 对于 $\forall t \in T$ 及 $P \subseteq C$,可定义 CL_t^{\geq} 与 CL_t^{\leq} 的分类系数为

$$\begin{cases} \alpha_p(CL_t^{\geq}) = \text{card}(\underline{PCL}_t^{\geq})/\text{card}(U), \\ \alpha_p(CL_t^{\leq}) = \text{card}(\underline{PCL}_t^{\leq})/\text{card}(U). \end{cases}$$

知识 P 关于划分 CL 的分类系数为: $\alpha_p(CL) = \sum_{t \in T} \{\alpha_p(CL_t^{\geq}) + \alpha_p(CL_t^{\leq})\}$ 。

参数2 对于 $t \in T$ 及 $\forall P \subseteq C$,可定义 CL_t^{\geq} 与 CL_t^{\leq} 的近似精度为

$$\begin{cases} \partial_p(CL_t^{\geq}) = \text{card}(\underline{PCL}_t^{\geq})/\text{card}(\overline{PCL}_t^{\geq}), \\ \partial_p(CL_t^{\leq}) = \text{card}(\underline{PCL}_t^{\leq})/\text{card}(\overline{PCL}_t^{\leq}). \end{cases}$$

知识 P 关于划分 CL 的近似精度为: $\partial_p(CL) =$

$\sum_{t \in T} \{\partial_p(CL_t^{\geq}) + \partial_p(CL_t^{\leq})\}$ 对于近似分类的不精确性, 可定义针对 P 而言 CL 划分的近似质量简称为分类质量为

参数3 知识 P 关于划分 CL 的分类质量为

$$\gamma_p(CL) = \frac{\sum_{t \in T} \text{card}(\underline{P}CL_t^{\geq}) + \sum_{t \in T} \underline{P}CL_t^{\leq})}{\sum_{t \in T} \text{card}(\overline{P}CL_t^{\geq}) + \sum_{t \in T} \overline{P}CL_t^{\leq}}.$$

它表示了决策表中所有 P -正确分类对象与所有 P -可能正确分类对象之比. 为了表示分类的相对近似, 可定义分类正域.

参数4 知识 P 关于划分 CL 的分类正域为

$$\beta_p(CL) = \frac{\text{card}(POS_p(CL))}{\text{card}(U)} = \frac{\sum_{t \in T} \text{card}(\bigcup_{X \in CL_t^{\geq}} \underline{P}X) + \sum_{t \in T} \text{card}(\bigcup_{X \in CL_t^{\leq}} \underline{P}X)}{\text{card}(U)}.$$

每一最小子集 $P \subseteq C$ 若满足 $\gamma_p(CL) = \gamma_c(CL)$ 则称为 CL 的一个约简, 记为 RED_{CL} . 即约简前后划分 CL 关于整个条件属性集合的分类质量保持不变. 一个带有准则的决策表不止一个约简, 所有约简的交称为该决策表的核. 记为 $CORE_{CL}$, $CORE_{CL} = \bigcap RED_{CL}$.

3.3 基于带有准则决策系统的数据约简算法(Data reduction algorithm based on decision system with criteria)

3.3.1 参数描述(Parameter description)

设决策系统中的条件属性 C 中有 m 个属性: C_1, C_2, \dots, C_m , 其中 m_1 个为有序属性, m_2 个为无序属性. 决策属性集合 D 的划分为 $CL = \{Cl_t, t \in T\}, T = \{1, 2, \dots, n\}$. 对每一个条件属性 C_i 计算以下 4 个参数: $\alpha_{C_i}(CL), \partial_{C_i}(CL), \gamma_{C_i}(CL), \beta_{C_i}(CL)$, 令 H_{C_i} 和 J_{C_i} 分别为这 4 个参数的算数平均和几何平均, $S_{C_i} = H_{C_i} + J_{C_i}$ 为条件属性的重要性系数.

3.3.2 算法过程(Algorithm procedure)

Step 1 输入决策表的条件属性与决策属性, 明确条件属性中的有序属性及无序属性并区分决策类;

Step 2 分别计算各个决策类关于整个条件属性的上下近似, 边界域正域和负域;

Step 3 对于条件属性集合 C , 计算 $V_c(CL)$;

Step 4 令初始约简属性集 C_0 为空;

Step 5 分别计算各个条件属性的参数, 计算其重要性系数 S_{C_i} , 并构成集合 S ;

Step 6 从 S 中取出最大值, 将对应属性加入约

简属性集 C_0 中, 计算 $V_{C_0}(CL)$, 并将 S_{C_i} 从 S 中删除;

Step 7 如果 $V_{C_0}(CL) < V_c(CL)$, 则转 Step 6, 否则转 Step 8;

Step 8 C_0 中包含的条件属性即为条件属性集合 C 的一个约简.

当不考虑属性的有序特征时, 此算法便可蜕变为标准粗糙集合的属性约简算法, 从而保证了带有有序属性的粗糙集数据约简与标准粗糙集的属性约简算法的一致性.

4 算例(Example)

在表 1 中, 电信网络质量评价指标一般由以下 6 个属性描述: $C1, C2, C3, C4, C5, D$ 分别表示网络负载能力、业务范围、信息传输时延、中心局址、拥塞标志、电信网络质量.

表 1 电信网络质量决策

Table 1 Telecommunication network quality decision

论域 U	$C1$	$C2$	$C3$	$C4$	$C5$	D
X1	高	多	长	L1	严重	好
X2	中	多	长	L1	严重	差
X3	中	多	长	L1	严重	好
X4	低	一般	长	L1	不能呼叫	差
X5	中	一般	短	L1	不能呼叫	差
X6	高	一般	短	L1	中等	好
X7	中	一般	长	L1	中等	好
X8	高	多	长	L2	中等	好
X9	中	多	长	L2	严重	好
X10	低	一般	长	L2	不能呼叫	差
X11	中	一般	短	L2	中等	好
X12	高	一般	短	L2	中等	好

由表可知, $C1, C2, C3, C5$ 是准则, 因为它的属性值是有序的, $C4$ 是一个属性, 其属性值范畴是无序的, D 是决策属性, 定义了两个有序决策, 详细地讲: 高 > 中 > 低, 多 > 一般, 短 > 长, 中等拥塞 > 严重拥塞 > 不能处理呼叫, 好 > 差, 其中 > 为“优于”符号.

4.1 基本粗集方法的约简结果(Reduction result based on rough set method)

由基本粗集方法中的基于正域算法, 定义 cl_1 为电信网络质量差决策类, cl_2 为电信网络质量好决策类. 经过计算, 决策规则为(计算过程略)

R1: ($C1=高 \rightarrow (\text{电信网络质量}=好})$;

R2: ($C1=中 \wedge (C4=L2 \rightarrow (\text{电信网络质量}=好})$);

R3: ($(C1=中 \wedge (C2=一般 \rightarrow (\text{电信网络质量}=好})$);

R4: ($(C1=中 \wedge (C2=多 \wedge (C4=L1 \rightarrow (\text{电信网络质量}=好}) \vee (\text{电信网络质量}=差})$);

R5: (C1=低)→(电信网络质量=差);

R6: (C5=不能处理呼叫)→(电信网络质量=差);

R7: (C1=中)∧(C3=短)∧(C4=L1)→(电信网络质量=差).

4.2 基于优先及不可分辨关系的近似约简结果(Approximation reduction result based on priority and indiscernibility relation)

定义决策类 cl_1^{\leq} 为电信网络质量最多为差的情况, cl_2^{\geq} 表示电信网络质量至少为好的情况, 因为只有两个决策类, 因而 $cl_1^{\leq}=cl_1$ 及 $cl_2^{\geq}=cl_2$. 计算 cl_1^{\leq} 及 cl_2^{\geq} 的 C —上下近似及边界域分别为

$$\underline{C}cl_1^{\leq}=\{4, 10\}, \overline{C}cl_1^{\leq}=\{2, 3, 4, 5, 7, 10\},$$

$$Bn_C(cl_1^{\leq})=\{2, 3, 5, 7\},$$

$$\underline{C}cl_2^{\geq}=\{1, 6, 8, 9, 11, 12\},$$

$$\overline{C}cl_2^{\geq}=\{1, 2, 3, 5, 6, 7, 8, 9, 11, 12\},$$

$$Bn_C(cl_2^{\geq})=\{2, 3, 5, 7\}.$$

cl_1^{\leq} 的近似精度为 0.63, cl_2^{\geq} 的近似精度为 0.86, 分类质量为 0.67. 只有一个约简即核, $Red_{cl}(c)=CORE_{cl}(c)=C_1, C_4, C_5$; 则最小决策规则为

R1: (C1=高)→(电信网络质量至少为好);

R2: (C1至少为中)∧(C4=L2)→(电信网络质量至少为好);

R3: (C1=低)→(电信网络质量至多为差);

R4: (C1=中)∧(C4=L1)→(电信网络质量=好)

∨(电信网络质量=差);

R5: (C5=不能处理呼叫)→(电信网络质量=差).

基于优先及不可分辨关系的粗集方法与基于不可分辨关系的粗集方法相比:

- * 规则数目减少, 规则利用较少的属性数及描述符, 更精炼;

- * 属性的约简数目减少;

- * 算法涵盖了属性了有序特征. 当不考虑属性的有序特征时, 此算法便可蜕变为标准粗糙集合的属性约简算法, 从而保证了带有有序属性的粗糙集数据约简与标准粗糙集的属性约简算法的一致性.

粗糙集合理论已经广泛应用于数据分析、图像处理、声音识别、决策支持分析和数据挖掘等领域, 并取得了很大成功. 本文在基本粗糙集合理论的基础上, 利用基于有序属性的数据挖掘方法, 能够快

速地在决策系统中挖掘出规则, 且挖掘出的规则简练, 更具合理性和可靠性.

参考文献(References):

- [1] PAWLAK Z. Rough set approach to knowledge-based decision support[J]. *European Journal of OR*, 1999, 3(27): 48–57.
- [2] PAWLAK Z. Rough sets[J]. *Int J of Computer and Information Science*, 1982, 11(5): 341–356.
- [3] ZIARKO W. Data-based acquisition and incremental modification of classification rules[J]. *Computer and Intelligence*, 1995, 4(11): 357–370.
- [4] 叶东毅, 黄翠微, 赵斌. 基于逼近精度的粗糙集属性约简算法[J]. 福州大学学报, 2000, 28(1): 8–10.
(YE Dongyi, HUANG Cuiwei, ZHAO Bin. An algorithm for attributions reduction in rough set based on approximation quality[J]. *Fuzhou University Journal*, 2000, 28(1): 8–10.)
- [5] ZIARKO W. Rough sets as a methodology for data mining[C]//POLKOWSKI L, SKOWRON A. *Rough sets in Knowledge Discovery: Methods Application*. Heidelberg: Physica-Verlag, 1998: 289–298.
- [6] 王珏. 粗糙集理论和统计学习理论[M]//陆然峰. 知识科学与计算科学. 北京: 清华大学出版社, 2003: 49–51.
(WANG Yu. Rough set theory and statistical learning theory[M]//LU Ranzheng. *Knowledge Science and Computing Science*. Beijing: Tsinghua University Press, 2003: 49–51.)
- [7] 张文修. 信息系统与知识发现[M]. 北京: 科学出版社, 2003.
(ZHANG Wenxing. *Information System and Knowledge Discovery*[M]. Beijing: Science Press, 2003.)
- [8] MIDEOLFART N, KOMOROWSKI J. A rough set approach to inductive logic programming[C]//ZIARKO W, YAO Y. *Rough Sets and Current Trends in Computing - Second International Conference*. Banff, Canada: Springer Press, 2000: 190–198.
- [9] MAHESWARA A, UMA V, SIROMONEY A, et al. The variable precision rough set inductive logic programming model and web usage graphs[J]. *New Frontiers in Artificial Intelligence-Joint JSAI 2001 Workshop Post-Proceedings*, 2001, 2253: 339–343.
- [10] BAZAN J G, SZCZUKA M. A collection of tools for rough set computations[C]//Proc of the 2nd Int Conf Rough Sets and Current Trends in Computing. Banff, Canada: Springer Press, 2000: 74–81.
- [11] ZHONG N, DONG J Z, OHSUGA S. Using background knowledge as a basis to control the rule discovery process[C]//ZIGHED D A, KOMOROWSKI J, ZYTKOW J. *Principles of Data Mining and Knowledge Discovery*. Berlin: Springer, 2001: 691–698.
- [12] GRECO S, MATARAZZO B, SLOWINSKI R. A new rough set approach to multicriteria and multiattribute classification[C]//Rough Set and Current Trends in Computing. Berlin: Springer, 1999: 57–61.

作者简介:

张文字 (1973—), 女, 博士, 西安邮电学院管理系副教授, 发表科技论文近30篇, 出版专著2本, 目前研究方向为数据挖掘、智能决策与网络可靠性, E-mail: zwy888459@sina.com 或 zwy888459@gmail.com.