

一种不一致不完备信息系统的最优选择及规则约简方法研究<sup>\*</sup>

伏明兰 曾黄麟

(四川理工学院 自贡 643000)

**摘要** 在分析不一致不完备信息系统规则提取的基础上,提出了先将不完备信息系统分为一致的和不一致的信息系统后再求其最优选择的方法。然后利用改进的分辨矩阵对所求得的不一致最优选择进行决策规则提取,并给出不确定决策规则的决策精度。最后,通过应用例子证明了本文提出方法的有效性。

**关键词** 粗糙集,不一致不完备信息系统,辨识矩阵,决策规则

## Optimization Selection and Rules Extraction in Inconsistent and Incomplete Information System

FU Ming-Lan ZENG Huang-Lin

(Si Chuan University of Science and Technology, Zigong 643000)

**Abstract** This paper analyzes the rules extraction in inconsistent and incomplete information system, and presents an advanced algorithm to acquire the optimized selection of incomplete information system which is divided into consistent subsystem and inconsistent subsystem. At the same time, a method based on an amended discernible matrix for rules extraction in inconsistent optimized selection is proposed. The possible rules precision in inconsistent information system is given. The effectiveness of the proposed method in the paper is demonstrated by an application example.

**Keywords** Rough sets, Inconsistent incomplete information system, Discernible matrix, Decision rules

## 1 引言

粗糙理论<sup>[1-3]</sup>于1982年由波兰学者 Z. Pawlak 提出后,已在机器学习,数据挖掘等若干领域得到了广泛应用。利用粗糙理论,能从信息系统中提取出确定性决策规则和非确定性决策规则。规则提取首先要求知识属性约简(知识约简),对于一致的信息系统,Skowron 提出的分辨矩阵是基于粗集求核和属性约简的主要方法之一。近年来,有不少学者以 Skowron 的分辨矩阵为基础定义新的分辨矩阵<sup>[4-8]</sup>,计算信息系统的核和属性约简,进行决策规则的提取,知识的挖掘。

对于不一致的不完备信息系统,由于其信息不仅不完备,而且不一致,但又广泛存在于实际生活中,因此,其属性约简和规则提取是粗糙集理论中的一个重要研究方向。近年来,人们在不一致不完备信息系统方面作了一些研究工作<sup>[1,9,10]</sup>,但尚存在一些问题未能解决,如算法占用内存空间大,时间复杂度高,有待于进一步改进。而且利用分辨矩阵对不一致信息系统进行规则提取时,只能求出确定性决策规则,不能求出带精度的不确定决策规则。

本文在分析这些问题的基础上,提出一种在时间和空间复杂度上都有较大改善的求取最优选择算法。还通过对分辨矩阵的改进,提出对于确定性、不确定性决策规则提取的方法,该方法能给出不确定决策规则的决策精度。

## 2 基本概念

**定义 1(决策信息系统)** 称  $S=(U, A, V, f)$  是一个决策信息系统。其中:  $U$  是非空有限论域,  $A=C \cup D$  是属性集,  $C$  和  $D$  分别是条件属性集和决策属性集。  $V=\bigcup_{a \in A} (V_a)$  是属性

值的集合,  $V_a$  表示属性  $a \in A$  的值域,  $f: U \times A \rightarrow V$  是一个信息函数。任意的对象  $x \in U$ , 对于任意属性  $a \in A$ , 对象  $x$  在属性  $a$  上的取值为  $a(x) \in V_a$ 。

**定义 2(不完备信息系统)** 在  $S=(U, A, V, f)$  中, 如果任意  $a(x)$  是唯一确定的, 则称  $S$  是完备的信息系统, 否则称  $S$  是不完备的信息系统。

**定义 3<sup>[1]</sup>(最大分布约简集)** 设  $S=(U, S, V, f)$  是不完备决策信息系统,  $B \subseteq A$ ,  $d$  为决策属性, 记

$$U/R_{\{d\}} = \{D_1, \dots, D_r\} \quad m_B(x) = \max\{D(D_j/S_B(x)) \mid j \leq r\} (x \in U),$$

其中  $D(E/F) = \frac{|E \cap F|}{|F|}$  是  $P(U)$  上的包含度,  $S_B(x) = \{y \in U \mid (x, y) \in SIM(B)\}$ ,  $S_B(x)$  表示  $x$  的相似类,  $SIM(B) = \{(x, y) \in U \times U \mid a_i(x) \cap a_i(y) \neq \text{null}, (a_i \in B)\}$ ,  $SIM(B)$  表示关于  $B$  的相似关系。

最大决策函数  $\gamma_B$  为

$$\gamma_B(x) = \{D_j \mid (D_j/S_B(x)) = m_B(x)\} (x \in U).$$

若对于任意  $x \in U$  有  $\gamma_B(x) = \gamma_A(x)$  成立, 则称  $B$  是  $(U, A, F, d)$  的最大分布协调集。若  $B$  是最大分布协调集, 且  $B$  的任何真子集都不是  $(U, A, F, d)$  的最大分布协调集, 称  $B$  是  $(U, A, F, d)$  的最大分布约简集。

**定义 4<sup>[1]</sup>(最优选择)** 设  $S=(U, A, V, f)$  是不完备决策信息系统,  $S^f=(U, A, f, d)$  是  $S$  的一个选择, 且  $B_f$  是  $S^f$  的最大分布约简集。若  $B_f$  是  $S$  的所有选择中的最小集合, 且满足  $\min_{x \in U} m_{B_f}(x) = \max \min_{x \in U} m_B(x)$ , 称缩减的完备信息系统  $(U, B_f, f, d)$  是  $S$  的最优选择。

**定义 5(对象之间的一致性定义)** 若决策表中的两个对

<sup>\*</sup> 基金项目: 四川省教育厅基础应用研究课题(2005A140)基金的部分资助。伏明兰 研究生, 主要研究领域为粗糙理论, 智能计算。曾黄麟 教授, 博导, 主要研究领域为粗糙理论, 智能计算。

象,满足如下条件:(1)条件属性的取值至少有一个不同;(2)有相同的条件属性取值时,决策属性的取值是相同的,称作这两个对象是一致的;否则称作不一致的。若决策表中任何一对对象都是一致的,称该决策表是一致的;否则称该决策表是不一致的。

**定义 6(分辨矩阵)** 设  $S=(U,A,V,f)$  是一个知识表达系统,研究论域  $U=\{x_1,x_2,\dots,x_m\}$ ,  $x_i(1\leq i\leq m)$  是研究对象,设属性集合  $A=\{a_1,\dots,a_n,d\}$ , 其中  $C=\{a_1,\dots,a_n\}$  为条件属性集,  $D=\{d\}$  决策属性集。对应的分辨矩阵  $M(C,D)$  是一个  $m\times m$  矩阵,其中  $m_{ij}(C,D)$  为分辨矩阵的第  $i$  行第  $j$  列元素。 $m_{ij}(C,D)\subseteq A, m_{ij}(C,D)$  的定义为:

$$m_{ij}(C,D)=\{a_k | a_k \in A; a_k(x_i) \neq a_k(x_j)\} i, j=1, 2, \dots, m; k=1, \dots, n+1$$

分辨矩阵的性质:

**性质 1** 设  $[x_i]_C$  表示  $x_i$  的  $C$ -等价类,在分辨矩阵中,  $x_i$  所在的行中的值为空 null 或为  $\{d\}$  的项所对应的列的对象所组成的对象集合  $X$  为  $x_i$  的  $C$ -等价类,记作  $X=[x_i]_C$ 。

证明:在分辨矩阵中,设  $x_i$  所在的行中值为空 null 或为  $\{d\}$  的项所在的列所对应的对象集为  $X$ ,根据分辨矩阵的定义,  $x_i$  与  $X$  中的元素不可分辨,即其条件属性的值都相同。反之,  $x_i$  若与某一元素  $x_j$  不可分辨,则它们在分辨矩阵中所对应的矩阵元素  $m_{ij}(C,D)$  为 null (决策属性相同) 或为  $\{d\}$  (决策属性不同)。综上所述,得  $X$  中的对象不可分辨,  $X$  为  $x_i$  的  $C$ -等价类。

**性质 2** 若分辨矩阵存在某个元素  $m_{ij}(C,D)=\{d\}$ , 则由  $x_i$  的  $C$ -等价类中的对象所得出的决策规则为不确定决策规则。

证明:由分辨矩阵的定义,若  $m_{ij}(C,D)=\{d\}$ , 则对象  $x_i$  和对象  $x_j$  的条件属性相同,而决策属性不同。根据定义 5,  $x_i$  和  $x_j$  是不一致的,因此由它得出的规则为不确定决策规则。

**性质 3** 若对象  $x_i$  属于正域  $POS_C(D)$ , 当且仅当在分辨矩阵  $M(C,D)$  中,  $x_i$  所在行的所有元素不为  $D=\{d\}$ 。

证明: $\Rightarrow$  设  $[x_i]_C, [x_i]_D$  分别表示在论域  $U$  中元素  $x_i$  的  $C$ -等价类,  $D$ -等价类,若  $x_i \in POS_C(D)$ , 则  $[x_i]_C \subseteq [x_i]_D$ , 由此得对于任意对象  $x_j$ , 若  $x_i$  与  $x_j$  的条件属性相同,则它们的决策属性肯定相同。根据定义 5, 对象  $x_i$  和  $x_j$  为一致的。而由分辨矩阵的定义,只有当出现不一致时,分辨矩阵的元素才为  $D=\{d\}$ , 由此得对象  $x_i$  所在行的所有元素不为  $D=\{d\}$ 。

$\Leftarrow$  若分辨矩阵中  $x_i$  所在的行的元素都不为  $D=\{d\}$ , 由分辨矩阵的定义 6, 只可能出现三种情况:条件属性和决策属性都相同;条件属性不同而决策属性相同;条件属性和决策属性都不相同,则根据定义 5,  $x_i$  与任意对象  $x_j$  都是一致的,从而得  $[x_i]_C \subseteq [x_i]_D$ , 由此得  $x_i \in POS_C(D)$ 。

**性质 4** 对于一个确定性或非确定性决策规则,在保持其条件属性  $C$ -等价类不变的前提下,其规则约简中所包含的属性与分辨矩阵中该规则对应的行的值不为 null 且不为  $\{d\}$  的每一项的交都不为空。

证明:对于一个决策规则  $x_i$ , 假设其规则约简  $redr$  中所包含的条件属性与分辨矩阵中  $x_i$  所在行的值不为 null 或  $\{d\}$  的某一项  $m_{ij}(C,D)$  的交为空, 则用  $redr$  中的属性表示该规则时,  $x_i$  与  $x_j$  由原来的可分辨变为不可分辨,从而改变了条件属性  $C$ -等价类, 违背保持条件属性  $C$ -等价类不变的前提。由反证法得证性质成立。

### 3 不一致不完备信息系统的最优选择及规则提取算法

文[1]中提出的求取不完备信息系统决策规则的方法为先找出不完备信息系统的最佳选择,再对最佳选择进行规则提取。最佳选择的目的是使信息表尽可能地趋于协调,也即选择使决策尽可能发生的条件。在此对文[1]中的方法作了改进,即先将信息系统分为一致信息系统和不一致信息系统,这样能大大减小算法的时间和空间复杂度。然后利用分辨矩阵进行规则提取。算法描述如下。

#### 3.1 不一致的不完备信息系统的最佳选择算法

(1)对不一致决策表进行划分,分为一致决策表 I 和非一致决策表 II;

(2)列出非一致决策表 II 的所有选择;

(3)对于第  $f$  个选择  $S^f, f=1, \dots, s, s$  为所有选择的个数。计算

$$m\ell(x)=\max\{D(D_j/[x]_C^f) | j \leq r\} (x \in U), r \text{ 为决策等价类的个数};$$

$$[x]_C^f = \{x_i | a_j(x) \cap a_j(x_i) \neq \text{null}, \forall a_j \in A, i=1, \dots, 9, j=1, \dots, 4\}, x \in U$$

(4)找出  $m = \max \min_{x \in U} m\ell(x)$ ;

(5)找到某个选择  $S^*$  使其满足  $\min_{x \in U} m\ell(x) = m$ ;

(6)将选择  $S^*$  与一致决策表 I 进行合并;

(7)合并后的选择可能有多个,求各个合并后的选择的最大分布约简集,并在其中挑选所含属性个数最小者;于是这个缩减的完备决策信息系统  $S'=(U, B_g, V, f)$  就是  $S=(U, A, V, f)$  的最佳选择,  $B_g$  为  $S$  的最大分布约简集。

#### 3.2 不一致的不完备最佳选择信息系统的决策规则提取算法

输入:经过属性约简后的不一致最佳选择决策表  $S'=(U, B_g, V, f)$ ;  
输出:约简后的决策规则,包括确定的决策规则和带精度的不确定决策规则

```

Begin;
根据定义 6, 求出不一致最佳选择  $S'$  的分辨矩阵。
求出决策等价类  $D_i, i=1, \dots, r$ 。
For  $i=1$  to  $m, \%m=|U|$ 
  If  $x_i$  未作判断
    If  $x_i$  所在分辨矩阵的行不包含  $\{d\}$  项
      ① 设该行中值为空 null 的项所在列对应的对象所组成的集合为  $X, X$  表示  $x_i$  的  $C$ -等价类。
      ② 在此行中选择属性个数最少的属性集, 使它与此行中的每一项的交不为空, 此属性值所对应的规则即为确定的最简决策规则。(可参考文[8]中的方法)
      ③ 输出确定性决策规则。
    else
      ① 扫描分辨矩阵中  $x_i$  所在行, 该行中值为 null 或为  $\{d\}$  的项所在列所对应的对象集合, 设为  $X, X$  表示  $x_i$  的  $C$ -等价类。
      ② 将  $X$  中的对象根据其决策值进行划分:  $X_1 \cup X_2 \cup \dots \cup X_l = X$ 。找出  $X_{\max}$  使  $|X_{\max}| = \max |X_j|, j=1, 2, \dots, l$ 。其中  $|X|$  表示集合  $X$  的基数。
      ③ 求出由集合  $X_{\max}$  中的任一对象  $x_m$  所得的决策规则, 并对该规则进行约简: 找出一最小属性集合  $Red$ , 使其与分辨矩阵中该行中除去值为 null 和  $\{d\}$  的项的交不为空。若某一项的值为  $\{a_i\}$  或  $\{a,d\}$ , 则  $a_i \in Red$ 。(注: 此步是在不改变原有条件属性  $C$ -等价类的前提下, 在分辨矩阵中即表现为不增加分辨矩阵中多余空集或多余  $\{d\}$  元素。)
      ④ 所得决策规则的精度为  $|X_{\max}|/|X|$ 。
    endif
  将  $X$  中的对象作上已判定的标记, 以后遇到  $X$  中的对象无需再进行规则提取。
endif
endif
endif
    
```

### 4 一个不一致不完备信息系统的最佳选择及其规则约简

一个不一致不完备决策信息系统  $S=(U, A, V, f)$  如表 1 所示。

表1 不一致不完备决策信息系统

U	a <sub>1</sub>	a <sub>2</sub>	a <sub>3</sub>	a <sub>4</sub>	d	m <sub>C</sub> <sup>f</sup> (x <sub>i</sub> )
x <sub>1</sub>	1	1	1	1	1	3/4
x <sub>2</sub>	2	1	1	1	1	1/2
x <sub>3</sub>	2	2	1,2	1	1	1
x <sub>4</sub>	2	1	2	1	2	1
x <sub>5</sub>	1,2	1	1,2	1	2	4/6
x <sub>6</sub>	1	1	1	1	2	3/4
x <sub>7</sub>	2	1	1,2	2	2	1
x <sub>8</sub>	1	1	1	1,2	2	3/4
x <sub>9</sub>	2	3	1	1	3	1

其中  $A=CUD, C=\{a_1, a_2, a_3, a_4\}$ , 为条件属性集,  $D=\{d\}$  为决策属性集,  $U/R_{(d)} = \{D_1 = \{x_1, x_2, x_3\}, D_2 = \{x_4, x_5, x_6, x_7, x_8\}, D_3 = \{x_9\}\}$

第①步: 求系统的最优选择

(1) 将不一致不完备决策信息系统分为一致决策表 I 和非一致决策表 II。

(2) 列出 II 的所有选择  $S^j$ , 共有  $2 \times 2 \times 2 = 8$  个选择, 由于篇幅限制, 在此不一一列举;

(3) 对于每个选择的每个对象  $x$  分别计算出  $m_C^f(x) = \max\{D(D_j/[x]_C^j) | j \leq r\}, r=1, 2, 3$ 。并找出每个选择中的  $\min_{x \in U} m_C^f(x)$  得:

- 选择  $S^1$ : 3/4      选择  $S^2$ : 2/3
- 选择  $S^3$ : 2/3      选择  $S^4$ : 1/2
- 选择  $S^5$ : 1/2      选择  $S^6$ : 1/2
- 选择  $S^7$ : 2/3      选择  $S^8$ : 1/2

(4)  $m = \max \min_{x \in U} m_C^f(x) = 3/4$ , 故得选择  $S^1$  为非一致决策表 II 的最优属性选择。

(5) 将选择  $S^1$  与一致决策表 I 合并后得到决策表如表 2 所示。

表2  $S^1$  与 I 合并后得到的决策表

U	a <sub>1</sub>	a <sub>2</sub>	a <sub>3</sub>	a <sub>4</sub>	d	m <sub>C</sub> <sup>f</sup> (x <sub>i</sub> )
x <sub>1</sub>	1	1	1	1	1	3/4
x <sub>2</sub>	2	1	1	1	1	1
x <sub>3</sub>	2	2	1,2	1	1	1
x <sub>4</sub>	2	1	2	1	2	1
x <sub>5</sub>	1	1	1	1	2	3/4
x <sub>6</sub>	1	1	1	1	2	3/4
x <sub>7</sub>	2	1	1,2	2	2	1
x <sub>8</sub>	1	1	1	1	2	3/4
x <sub>9</sub>	2	3	1	1	3	1

此决策表有 4 个选择, 设为  $S_{11}, S_{12}, S_{13}$  和  $S_{14}$ 。

(6) 分别求出它们的最大分布约简集:

- $S_{11} : \{a_1, a_2, a_3, a_4\}$      $S_{12} : \{a_1, a_2, a_3\}$
- $S_{13} : \{a_1, a_2, a_3, a_4\}$      $S_{14} : \{a_1, a_2, a_3\}$

(7) 由此得最终的最优选择有 2 个, 在  $S_{12}$  和  $S_{14}$  中任选 1 个进行规则约简, 在此选择  $S_{12}$ , 如表 3 所示。

表3 选择  $S_{12}$

U	a <sub>1</sub>	a <sub>2</sub>	a <sub>3</sub>	d	m <sub>C</sub> <sup>f</sup> (x <sub>i</sub> )
x <sub>1</sub>	1	1	1	1	3/4
x <sub>2</sub>	2	1	1	1	1
x <sub>3</sub>	2	2	1	1	1
x <sub>4</sub>	2	1	2	2	1
x <sub>5</sub>	1	1	1	2	3/4
x <sub>6</sub>	1	1	1	2	3/4
x <sub>7</sub>	2	1	2	2	1
x <sub>8</sub>	1	1	1	2	3/4
x <sub>9</sub>	2	3	1	3	1

第 2 步: 由分辨矩阵求规则约简求出表 3 所示的最优选择不一致信息系统的分辨矩阵, 如表 4 所示。

表4 不一致最优选择的分辨矩阵

U	x <sub>1</sub>	x <sub>2</sub>	x <sub>3</sub>	x <sub>4</sub>	x <sub>5</sub>	x <sub>6</sub>	x <sub>7</sub>	x <sub>8</sub>	x <sub>9</sub>
x <sub>1</sub>	null	a <sub>1</sub>	a <sub>1</sub> a <sub>2</sub>	a <sub>1</sub> a <sub>2</sub> d	d	d	a <sub>1</sub> a <sub>3</sub> d	d	a <sub>1</sub> a <sub>2</sub> d
x <sub>2</sub>	a <sub>1</sub>	null	a <sub>2</sub>	a <sub>3</sub> d	a <sub>1</sub> d	a <sub>1</sub> d	a <sub>3</sub> d	a <sub>1</sub> d	a <sub>2</sub> d
x <sub>3</sub>	a <sub>1</sub> a <sub>2</sub>	a <sub>2</sub>	null	a <sub>2</sub> a <sub>3</sub> d	a <sub>1</sub> a <sub>3</sub>	a <sub>1</sub> a <sub>2</sub> d	a <sub>2</sub> a <sub>3</sub> d	a <sub>1</sub> a <sub>2</sub> d	a <sub>2</sub> d
x <sub>4</sub>	a <sub>1</sub> a <sub>3</sub> d	a <sub>3</sub> d	a <sub>2</sub> a <sub>3</sub> d	null	a <sub>1</sub> a <sub>3</sub>	a <sub>1</sub> a <sub>3</sub>	null	a <sub>1</sub> a <sub>3</sub>	a <sub>2</sub> a <sub>3</sub> d
x <sub>5</sub>	d	a <sub>1</sub> d	a <sub>1</sub> a <sub>3</sub>	a <sub>1</sub> a <sub>3</sub>	null	null	a <sub>1</sub> a <sub>3</sub>	null	a <sub>1</sub> a <sub>2</sub> d
x <sub>6</sub>	d	a <sub>1</sub> d	a <sub>1</sub> a <sub>2</sub> d	a <sub>1</sub> a <sub>3</sub>	null	null	a <sub>1</sub> a <sub>3</sub>	null	a <sub>1</sub> a <sub>2</sub> d
x <sub>7</sub>	a <sub>1</sub> a <sub>3</sub> d	a <sub>3</sub> d	a <sub>2</sub> a <sub>3</sub> d	null	a <sub>1</sub> a <sub>3</sub>	a <sub>1</sub> a <sub>3</sub>	null	a <sub>1</sub> a <sub>3</sub>	a <sub>2</sub> a <sub>3</sub> d
x <sub>8</sub>	d	a <sub>1</sub> d	a <sub>1</sub> a <sub>2</sub> d	a <sub>1</sub> a <sub>3</sub>	null	null	a <sub>1</sub> a <sub>3</sub>	null	a <sub>1</sub> a <sub>2</sub> d
x <sub>9</sub>	a <sub>1</sub> a <sub>2</sub> d	a <sub>2</sub> d	a <sub>2</sub> d	a <sub>2</sub> a <sub>3</sub> d	a <sub>1</sub> a <sub>2</sub> d	a <sub>1</sub> a <sub>2</sub> d	a <sub>2</sub> a <sub>3</sub> d	a <sub>1</sub> a <sub>2</sub> d	null

对未作已判断标记的对象进行规则提取。对象  $x_1, x_5, x_6, x_8$  所在的行包含值为  $\{d\}$  的元素, 因而, 由分辨矩阵的性质 2 得由它们所得的决策规则为不确定决策规则; 对象  $x_2, x_3, x_4, x_7, x_9$  所在的行不包含值为  $\{d\}$  的元素, 因而能得确定的决策规则, 决策精度为 1。

对象  $x_1$  所在的行有 null 元素及  $\{d\}$  元素, 所对应的对象集为  $\{x_1, x_5, x_6, x_8\}$ , 由性质 1 得  $x_1$  的 C-等价类为:  $\{x_1, x_5, x_6, x_8\}, |\{x_1, x_5, x_6, x_8\}| = 4$ 。其中,  $\{x_1\} \subset D_1 (d=1), |\{x_1\}| = 1; |\{x_5, x_6, x_8\} \subset D_2 (d=2) | \{x_5, x_6, x_8\}| = 3 > |\{x_1\}| = 1$ , 故得带精度的决策规则:  $(a_1=1) \wedge (a_2=1) \wedge (a_3=1) \Rightarrow d=2 (3/4)$ 。将对象  $\{x_1, x_5, x_6, x_8\}$  作上已判断的标记, 以后不需再对其进行规则提取。在不改变原有不协调条件属性等价类的前提下, 即不增加分辨矩阵第  $x_1, x_5, x_6, x_8$  行的多余空集或多余  $\{d\}$  元素的前提下, 可得最简决策规则:  $(a_1=1) \Rightarrow d=2 (3/4)$  ( $\{a_1\}$  与分辨矩阵第  $x_1, x_5, x_6, x_8$  行的 null 非或  $\{d\}$  的元素的交都不为空)。

对象  $x_2, x_3, x_4, x_7, x_9$  在分辨矩阵中所在的行不包含值为  $\{d\}$  的元素, 定理 1 得由它们所得的决策规则为确定的。规则约简的获得方法为: 在不改变原有 C-等价类的前提下, 在每一个确定性决策规则所对应的分辨矩阵的行中找出最小属性个数的属性集 X, 使 X 与该行的每一个非 null 或  $\{d\}$  的元素的交都不为空。X 所对应的属性集即为约简的决策规则, 可参考文[8]中提出的方法。其决策规则如下:

- $x_2 : (a_1=2) \wedge (a_2=1) \wedge (a_3=1) \Rightarrow d=1$
- $x_3 : [(a_1=2) \wedge (a_2=2)] \vee [(a_2=2) \wedge (a_3=1)] \Rightarrow d=1$
- $x_4, x_7 : a_3=2 \Rightarrow d=2$
- $x_9 : (a_2=3) \Rightarrow d=3$

结束语 本文基于概率统计概念, 将不一致不完备信息系统分为一致信息系统和不一致信息系统, 再找出非一致信息系统的最优选择, 并与一致决策表合并, 重建缺省数据, 从而进行经典粗集模型支持下求解不完备信息系统的知识处理。对于求不一致最优选择的规则约简必须在不改变原有条件属性 C-等价类的前提下进行, 此前提是一个较为严格的条件。若条件有变, 其相应的属性约简和规则提取方法也会在本文提出的方法的基础上作一些改变, 这都在我们的进一步研究当中。

参考文献

1 张文修, 仇国芳. 基于粗糙集的不确定决策(154-157)[M]. 北京: 清华大学出版社, 2005

大规模数据集下支持向量机训练样本的缩减策略<sup>\*</sup>)罗 瑜<sup>1</sup> 易文德<sup>2</sup> 王丹琛<sup>3</sup> 何大可<sup>1</sup>(西南交通大学信息科学与技术学院 成都 610031)<sup>1</sup> (重庆文理学院数学与计算机科学系 永川 402160)<sup>2</sup>  
(四川省信息安全测评中心 成都 610017)<sup>3</sup>

**摘要** 大量数据下支持向量机的训练算法是 SVM 研究的一个重要方向和焦点。该文从分析 SVM 训练问题的实质和难点出发,提出一种在训练前先求出类别质心,去除非支持向量对应的样本,从而达到缩小样本集的方法。该方法在不损失分类正确率的情况下具有更快的收敛速度,并从空间几何上解释了支持向量机的原理。仿真实验证明了该方法的可行性和有效性。

**关键词** 支持向量机,分解算法,类别质心,准支持向量

## Sample Reduction Strategy for Support Vector Machines with Large-Scale Data Set

LUO Yu<sup>1</sup> YI Wen-De<sup>2</sup> WANG Dan-Chen<sup>3</sup> HE Da-Ke<sup>1</sup>(School of Information Science & Technology, Southwest Jiaotong University, Chengdu 610031)<sup>1</sup>(Dept. of Mathematics & Computer Science, Chongqing University of Arts and Sciences, Yongchuan 402160)<sup>2</sup>(Sichuan Province Information Security Testing Evaluation Center, Chengdu 610017)<sup>3</sup>

**Abstract** Training algorithm for large-scale support vector machines(SVM) is an important and active subject in the field of SVM research. After the analysis of the nature and difficulties in training SVM, a new reduction strategy is proposed in this paper for training svm with large-scale sample set. In general, class centroid is solved before training and removing the samples corresponding to non support vectors. Through this method, the number of samples is reduced before training svm. This method is fast in convergence without accurate loss and propose the explanation of SVM theory from space geometry. The results of simulation experiments show the feasibility and effectiveness of this method.

**Keywords** Support vector machines, Decomposition algorithm, Reduction strategy, Centroid, Quasi-support vectors

## 1 引言

Vapnik<sup>[1,2]</sup>等人提出的统计学习理论是一种针对小样本的学习理论,它避免了人工神经网络等方法的网络结构难以确定、过学习和欠学习以及局部极小等问题,被认为是目前针对小样本的分类、回归等问题的最佳理论。

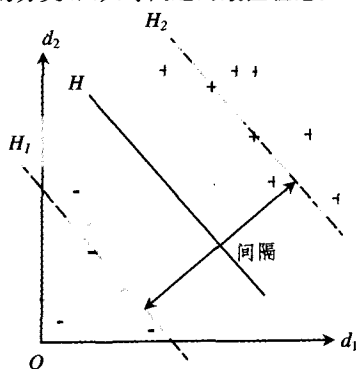


图1 分类超平面

支持向量机的基本思想是对于非线性可分样本,基于1909年 Mercer 核展开定理,将其输入样本通过非线性变换映射到另一个高维空间  $Z$  (Hilbert 空间)中,在变换后的空间中构造一个最优的分类超平面  $H$  (图1),使其在保证分类精度(满足经验风险)的同时最大化超平面两侧的空白区域(即最大化置信范围),也即是  $H_1$  与  $H_2$  间的几何间隔(图1)。这使得分类的结果不但在训练集上得到优化,而且在整个样本集上的风险也有上界,这就是 SVM 的结构风险最小化思想。鉴于篇幅的原因,本文只针对两类问题进行探讨,但本文的结论同样适用于多类问题。

## 2 支持向量机理论基础

对于两类问题,给定样本集  $(x_i, y_i); x_i \in R^n, y_i = \pm 1, i = 1, 2, \dots, l$  和核函数  $K(x_i, x_j)$ 。  $K$  对应特征空间  $Z$  (Hilbert 空间)的内积,  $K(x_i, x_j) = (\Phi(x_i) \cdot \Phi(x_j))$ 。设计基于 SVM 的分类器,就是在  $Z$  中寻找最优超平面  $H$ ,即最大化两类的几何间隔。对于样本集线性可分和不可分,都可用下式来表达<sup>[1,3]</sup>:

<sup>\*</sup>)基金项目:上海市特种光纤重点实验室科研项目,地铁 CBTC 无线接入安全认证算法研究。罗 瑜 博士生,研究方向为并行计算及模式识别。

- 曾黄麟. 粗集理论及其发展[M]. 重庆:重庆大学出版社,1998
- 曾黄麟. 智能计算[M]. 重庆:重庆大学出版社,2004
- 郑书富,管延勇,史开泉. 分辨矩阵与它在非一致决策中的应用[J]. 山东大学学报,2005,35(2):86~89
- Xiaohua H, Cercon N. Learning in relational data-bases: a rough set approach[J]. Computation Intelligence, 1995, 11(2): 323~337
- 叶东毅, 陈昭炯. 一个新的差别矩阵及其求核方法[J]. 电子学报, 2002, 31(7): 1086~1088
- 王希雷,王磊. 粗集中区分矩阵对不一致问题处理的研究[J]. 微机发展, 2003, 13(6): 119~120
- 汪廷华,程从从. Rough 集中基于分明矩阵的决策规则约简研究[J]. 计算机科学, 2005, 32(8, A): 42~44
- 黄兵,周献中. 不完备信息系统分配约简与规则提取的矩阵算法[J]. 计算机工程, 2005, 31(17): 20~22
- 管延勇,薛佩军,王洪凯. 不完备信息系统的可信决策规则提取与 E-相对约简[J]. 系统工程理论与实践, 2005, 12: 76~82