

# 不协调目标信息系统中基于改进差别信息树的分布属性约简

龙柄翰<sup>1</sup> 徐伟华<sup>2</sup> 张晓燕<sup>2</sup>

(重庆理工大学理学院 重庆 400054)<sup>1</sup> (西南大学数学与统计学院 重庆 400715)<sup>2</sup>

**摘 要** 在信息系统不协调的背景下,文中研究了如何有效地求解分布属性约简的问题。利用分布协调集的判定定理,提出了一种在不协调目标信息系统背景下进行分布属性约简的新方法。受到差别矩阵和差别信息树的启发,在该方法中构造了一种利用改进的差别信息树进行分布属性约简的算法。该信息树实现了对差别矩阵中的非空元素以及冗余信息的压缩储存,极大简化了时间复杂度及空间复杂度。

**关键词** 不协调信息系统,分布属性约简,分布协调集,改进差别信息树

中图法分类号 TP181 文献标识码 A

## Distribution Attribute Reduction Based on Improved Discernibility Information Tree in Inconsistent System

LONG Bing-han<sup>1</sup> XU Wei-hua<sup>2</sup> ZHANG Xiao-yan<sup>2</sup>

(School of Science, Chongqing University of Technology, Chongqing 400054, China)<sup>1</sup>

(School of Mathematics and Statistics, Southwest University, Chongqing 400715, China)<sup>2</sup>

**Abstract** Under the background of inconsistent systems, this paper studied how to effectively solve the problem of distributed attribute reduction. By using the judgment theorem of distributed coordination set, a new method of distributed attribute reduction under the background of inconsistent system was proposed. Inspired by difference matrix and discernibility information tree, in this method, an algorithm is constructed which uses the improved discernibility information tree to reduce the distribution attribute. The information tree realizes the compression and storage of non-empty elements and redundant information in the discernibility matrix, and greatly simplifies the time complexity and the space complexity.

**Keywords** Inconsistent system, Distribution attribute reduction, Distribution coordination set, Improved discernibility information tree

## 1 引言

粗糙集是一种进行数据分析的数学工具,它可以对不完整或者不确定的数据进行处理。特别是对一些不完备、不准确或者模糊的信息,它能挖掘出其中的隐藏信息,并发现其内在规律。粗糙集是波兰学者 Pawlak 最先提出的,他详尽、系统地阐述了粗糙集理论的相关概念,也为粗糙集思想的发展与推广打下了坚实的基础。如今,此理论已广泛应用于各个领域,如知识获取、人工智能、数据挖掘、过程控制等。由于粗糙集理论中最基础、最重要的研究课题就是属性约简,因此,属性约简受到了众多学者的广泛关注<sup>[1-3]</sup>。

属性约简,也称为知识约简或者特征选择。一般来说,知识库中的知识重要程度不是完全相同的,有些知识是必要的,有些知识是不必要的,还有些知识甚至是冗余的。正是由于知识的重要程度不同,才需要利用属性约简对知识进行压缩分类。属性约简是从条件属性集中去除不相关或不重要的冗余条件属性,同时保持知识库的分类能力不变。属性约简首

先需要从信息系统中提取有用的信息,然后简化知识,让存储空间变得更少并降低处理时间。经过属性约简后的信息系统更加简单清晰,方便决策者做出正确有效的决策。属性约简的目的是对原始数据进行降维处理,使存储空间大幅减少,并提高分类准确率与速度<sup>[4-5]</sup>。

关于属性约简,很多学者已经提出了不同的约简算法<sup>[6-8]</sup>。文献[6]在变精度的情况下,应用了属性约简的定义,在启发式算法的基础上加入了属性核的思想。文献[7]提出的算法效率低,步骤太繁琐。文献[8]提出的算法应用了属性约简的性质,并加入了一些启发性知识,但占用空间仍然较大。本文对各种算法进行对比研究,得到了一种利用改进差别信息树进行分布属性约简的算法<sup>[3,9]</sup>。

在处理实际应用数据的过程中,有时会遇到决策表的预处理数据值不一致的情况。很多原因都可以造成这种情况,例如:条件属性太单薄、测量样本属性值时有误差等。这就会使得信息系统变得不协调。因此,近年来许多学者都致力于不协调信息系统属性约简的研究<sup>[10-11]</sup>中。

本文受国家自然科学基金项目(61472463, 61402064, 61772002),重庆市自然科学基金项目(cstc2015jcyjA40053),重庆市教委科技项目(KJ1709221)资助。

龙柄翰 男,硕士生,主要研究方向为粗糙集理论;徐伟华 男,博士,教授,博士生导师,主要研究方向为人工智能与粒计算、模糊数学、信息科学, E-mail: datongxuwei@126.com(通信作者);张晓燕 女,博士,副教授,硕士生导师,主要研究方向为粒计算与人工智能、概念格与不确定性推理。

文献[8,12-13]在不协调信息系统背景下,提出了分配约简、最大分布约简、分布约简等概念,并详尽地介绍了它们的定义与性质。本文在保持分布协调集不变的基础上对不协调信息系统进行分布约简。

差别矩阵是属性约简的一种重要工具,它简单清晰,但是冗余元素多、占用的存储空间大、计算约简复杂度高<sup>[14]</sup>。文献[15]提出了一种树型结构,即差别信息树。该树是一棵基于条件属性从左到右序列的有序前缀树。在差别信息树中,差别矩阵中的相同元素被映射到同一路径上,并且部分具有子、父集关系的元素被映射到了子集所对应的路径上,若这些路径拥有共同的节点,则它们共享路径前缀,让共同重复的冗余节点信息占用一条路径,达到降低存储空间的目的。因此,相比差别矩阵,差别信息树的存储复杂度较低,并包含了构建属性约简所需要的所有差别信息<sup>[1,7,16]</sup>。

然而,差别信息树并没有对差别矩阵中核属性元素压缩存储。因此,本文提出了一种改进差别信息树的方法。该方法利用了核属性剪枝策略并对属性重要度进行降序排列。其优势在于:核属性剪枝策略可以确保改进的差别信息树能消除所有核属性的父集,进一步压缩存储空间;属性重要度降序策略确保了越重要的属性越靠近树的根节点,并使得更多的属性在改进的差别信息中共享前缀和路径。这样既压缩了差别矩阵中元素的存储空间,又可以确保启发式算法可以得出条件属性数量最少的约简<sup>[17-18]</sup>。

## 2 基本概念

为了后文的叙述方便,本节给出了与本文相关的基本概念和内容。

定义 1<sup>[5]</sup> 差别信息树是一棵有序树,树中每个节点至多只有 $|C|$ 棵子树( $C$ 为决策表条件属性集),而且必须满足以下条件。

1)差别信息树中每个节点主要包含以下4结构:属性名、父亲指针、孩子指针和同名指针。其中,该节点所对应的条件属性称为属性名;指向该节点的父亲节点的指针称为父亲指针;指向该节点的孩子节点的指针称为孩子指针;指向其他路径中与该节点具有相同属性名节点的指针称为同名指针。

2)属性指针头表由两个部分构成,分别是:属性名和同名指针。其中,标识属性指针头表所对应的条件属性称为属性名;指向差别信息树中与属性指针头表项具有相同属性名的最左边节点的指针叫做同名指针。

定义 2<sup>[3]</sup> 称 $\langle U, A, F, D, G \rangle$ 是决策信息系统,其中, $U$ 是样本集合, $U = \{x_1, x_2, \dots, x_n\}$ ;  $A$ 是条件属性集合, $A = \{a_1, a_2, \dots, a_m\}$ ;  $D$ 是决策属性集合, $D = \{d_1, d_2, \dots, d_q\}$ ;  $F$ 是 $U$ 与 $A$ 的函数关系集合。其中, $F = \{f_k: U \rightarrow V_k, k \leq m\}$ ,  $V_k$ 是样本 $x_k$ 在 $a_k$ 的值域, $G$ 是 $U$ 与 $D$ 的函数关系集合, $G = \{g_k: U \rightarrow V_k, k \leq q\}$ ,  $V_k$ 是样本 $x_k$ 在 $d_k$ 的值域。对于任意的 $B \subseteq A$ 的不可分辨关系:

$$R_B = \{(x, y) : f_k(x) = f_k(y), \forall a_k \in B\}$$

$$R_D = \{(x, y) : g_k(x) = g_k(y), \forall a_k \in D\}$$

它们在 $U$ 上产生的划分分别为:

$$U/R_B = \{[x]_B : x \in U\}$$

$$U/R_D = \{[x]_D : x \in U\} = \{D_1, D_2, \dots, D_r\}$$

其中:

$$[x]_B = \{y : (x, y) \in R_B\}$$

$$[x]_D = \{y : (x, y) \in R_D\}$$

定义 3<sup>[3]</sup> 给定决策信息系统 $\langle U, A, F, D, G \rangle$ ,如果有 $R_A \subseteq R_D$ ,则称该信息决策系统是协调的,否则称该信息决策系统是不协调的。

定义 4<sup>[3]</sup> 设 $\langle U, AT \cup \{d\}, F, G \rangle$ 为信息决策系统, $A \subseteq AT, x \in U$ ,记:

$$D(D_j/[x]_A) = \frac{|D_j \cap [x]_A|}{|[x]_A|}, j \leq r$$

则 $D$ 是 $P(U)$ 上的包含度,若记:

$$\mu_A(x) = (D(D_1/[x]_B), \dots, D(D_r/[x]_B))$$

则称 $\mu_A(x)$ 为对象 $x$ 关于属性集 $B$ 在目标信息系统中的广义决策分布函数。显然, $\mu_A(x)$ 是 $U/R_D$ 上的条件概率分布。

定义 5<sup>[19]</sup> 设 $I = \langle U, AT \cup \{d\}, F, G \rangle$ 为信息决策系统, $A \subseteq AT$ 。若对于任意 $x \in U$ ,都有 $\mu_A(x) = \mu_{AT}(x)$ ,则称 $A$ 是 $I$ 中关于 $R_{AT}$ 的分布协调集,简称分布协调集。进而,若 $A$ 的任何真子集都不是分布协调集,则称 $A$ 为 $I$ 中关于 $R_{AT}$ 的分布约简,简称分布约简。

定义 6<sup>[19]</sup> 设 $I = \langle U, AT \cup \{d\}, F, G \rangle$ 为信息决策系统。记:

$$Dis_{AT}^A(x_i, x_j) =$$

$$\begin{cases} \{a \in AT : f(x_i, a) \neq f(x_j, a)\}, & \mu_A(x_i) \neq \mu_A(x_j) \\ \phi, & \mu_A(x_i) = \mu_A(x_j) \end{cases}$$

则称 $Dis_{AT}^A(x_i, x_j)$ 为 $I$ 中对象 $x_i, x_j$ 关于关系 $R_{AT}$ 的分布相对差别属性集,简称分布差别属性集。记:

$$Dis_{AT}^A = (Dis_{AT}^A(x_i, x_j))_{|U \times U|}$$

称 $Dis_{AT}^A$ 为 $I$ 中关于关系 $R_{AT}$ 的分布相对差别矩阵,简称分布差别矩阵。

定义 7<sup>[3]</sup> 设 $I = \langle U, AT \cup \{d\}, F, G \rangle$ 为信息决策系统,对象关于关系 $R_{AT}$ 的分布差别矩阵为 $Dis_{AT}^A$ 。记:

$$M_{AT}^A = \bigwedge \{ \bigvee \{a \mid a \in Dis_{AT}^A(x_i, x_j)\} \mid \forall x_i, x_j \in U \}$$

则称 $M_{AT}^A$ 为 $I$ 关于关系 $R_{AT}$ 的分布差别公式。

定义 8<sup>[3]</sup> 设 $I = \langle U, AT \cup \{d\}, F, G \rangle$ 为信息决策系统。分布差别公式 $M_{AT}^A$ 的极小析取范式为:

$$M_{\min}^A = \bigvee_{k=1}^p \left( \bigwedge_{s=1}^{q_k} a_s \right)$$

记:

$$B_\mu^k = \{a_s \mid s=1, 2, 3, \dots, q_k\}$$

则 $\{B_\mu^k \mid k=1, 2, 3, \dots, p\}$ 是所有分布约简构成的集合。

定义 9 设 $\langle U, AT \cup \{d\}, F, G \rangle$ 为决策信息系统, $Dis_{AT}^A = (Dis_{AT}^A(x_i, x_j))_{|U \times U|}$ 为差别矩阵,属性 $a \in AT$ ,则 $a$ 的重要度为:  
 $SGF(a, R, D) = p(a)$

其中, $p(a)$ 为差别矩阵中,构建差别信息树所用到的所有差别信息中属性 $a$ 所出现的频率。

## 3 不协调信息系统中改进差别信息树的设计与实现

众所周知,信息系统决策表中的“核”是一个非常重要的条件属性集的子集,表示为 $Core_C(D)$ 。对于差别矩阵而言,差别矩阵中所有只包含一个元素的集合的“并”构成决策表的核。核之所以重要,是因为核中任意属性被删除都会影响决策表的分类能力。基于该分析,在压缩存储差别矩阵非空元素时,如果让决策表的核作为启发信息,那么在构建改进差别

信息树时,不仅不需要差别矩阵中晚于核属性出现的父集元素参与树的构建,而且能消除树中已构建且包含核属性的路径,从而,实现了消除差别矩阵中所有核属性父集冗余元素的目的。核属性剪枝策略虽然保证了能消除核属性的父集差别信息,但是不能确保其他父集差别信息能被删除。

属性重要度用来表示这些属性元素在决策表中的重要程度。因为,属性越重要,其在决策表中所处位置越高,反之亦然。如果在构建改进差别信息树时,加入属性重要度作为启发式信息,根据上面的分析可知:它不但能使越重要的属性越临近根节点,加大改进差别信息树的压缩程度,还能有效地消除父集差别信息。基于以上分析,本文利用属性重要度策略设计了基于核属性剪枝策略和属性重要度策略的改进信息差别信息树方法,如算法 1 所示。

**算法 1 基于核属性剪枝策略和属性重要度策略的改进型差别信息树构建算法**

输入:决策表 T

输出:改进型差别信息树(IDI<sub>tree</sub>)和 Core<sub>c</sub>(D)

```

{
1. 令 Corec(D) = ∅;
2. 计算条件属性集中属性重要度,并令 CList 为按属性重要度降序排列的序列;
3. 创建改进差别信息根节点 TN,并令 TN 为 null;
4. 对决策表 T 中每个对象对 ⟨xi, xj⟩, 计算其差别信息 DI, 并按序列 CList 排列 DI 中的条件属性, 表示为 order(DI);
   if(order(DI) ∩ Corec(D) = ∅)
   {
     if(|order(DI)| = 1) Corec(D) = Corec(D) ∪ order(DI)
     且根据属性指针头表删除含有该核属性的所有路径;
     while(order(DI) ≠ ∅)
     {
       4.1 令属性 b 是排列在 order(DI) 中最左边的元素;
       4.2 if(TN 有一孩子节点 N, 且 N 的属性名为 b)
       {
         如果 N 是一叶子节点, 返回; 否则 TN = N;
       }
       else
       {
         4.2.1 创建一新节点 N', 节点 N' 作为 TN 一子节点, 同时初始 N' 的属性名为 b, 并通过该节点的同名指针连接到具有与该节点有相同属性名的节点上, 从而构成了一个同名属性节点链;
         4.2.2 令 TN = N';
       }
       4.2.3 order(DI) = order(DI) - {b};
     }
   }
}

```

为了验证本文所提改进差别信息树的有效性,基于改进差别信息树,采用至下而上的策略,提出了一种基于改进差别信息树的粗糙集属性约简算法。算法具体描述如算法 2 所示。

**算法 2 基于改进差别信息树的属性约简算法**

输入:改进差别信息树

输出:决策表的一个属性约简

1. 创建一个空集 R。
2. 获取改进差别信息头表 HT, 并令 CList 为头表中属性至上而下的有序序列集合。

3. 令集合 Core<sub>c</sub>'(D) 是改进差别信息树中只含一个节点的路径所对应的属性所构成的集合。
4. 如果 Core<sub>c</sub>'(D) 非空, 则对于所有 a ∈ Core<sub>c</sub>'(D), 在改进差别信息树中删除包含 a 属性的所有路径。
5. 令 R ← Core<sub>c</sub>'(D) 且 CList ← CList - Core<sub>c</sub>'(D)。
6. While  
   CList ≠ ∅ ∧ IDI<sub>tree</sub> ≠ ∅  
   do  
   {  
     6.1 选取 CList 中最右边的一属性 R, 且令 CList ← CList - {C<sub>i</sub>};  
     6.2 如果 HT[C<sub>i</sub>] ≠ ∅ 并且 C<sub>i</sub> ∉ R, 则根据头表 HT[C<sub>i</sub>] 所对应的同名指针, 搜索同名指针构成的指针链, 在搜索过程中, 若链中节点为叶子节点, 则在改进差别信息树中删除该叶子节点;  
     6.3 把改进差别信息树中的只有一个节点的路径所对应的属性加入集合 R 中; 并在改进差别信息树中删除包含这些属性的路径; 同时在 CList 中删除这些属性。  
   }  
- 7. 输入 R, 算法结束。

**4 案例分析**

例 1 表 1 为一不协调决策表, 其中, 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>}, C = {a, b, c, d} 且 D = {e}。

表 1 不协调决策信息系统

| U              | a | b | c | d | e |
|----------------|---|---|---|---|---|
| x <sub>1</sub> | 2 | 1 | 3 | 1 | 1 |
| x <sub>2</sub> | 2 | 1 | 3 | 2 | 1 |
| x <sub>3</sub> | 2 | 2 | 2 | 1 | 1 |
| x <sub>4</sub> | 2 | 2 | 2 | 1 | 2 |
| x <sub>5</sub> | 2 | 2 | 2 | 2 | 2 |
| x <sub>6</sub> | 1 | 3 | 3 | 3 | 2 |

首先, 计算不协调信息系统中 x 关于属性集 C 的广义分布函数。[x<sub>1</sub>]<sub>C</sub> = {x<sub>1</sub>}, [x<sub>2</sub>]<sub>C</sub> = {x<sub>2</sub>}, [x<sub>3</sub>]<sub>C</sub> = {x<sub>3</sub>, x<sub>4</sub>}, [x<sub>5</sub>]<sub>C</sub> = {x<sub>5</sub>}, [x<sub>6</sub>]<sub>C</sub> = {x<sub>6</sub>}。D<sub>1</sub> = {x<sub>1</sub>, x<sub>2</sub>, x<sub>3</sub>}, D<sub>1</sub> = {x<sub>4</sub>, x<sub>5</sub>, x<sub>6</sub>}。μ<sub>A</sub>(x<sub>1</sub>) = (1, 0), μ<sub>A</sub>(x<sub>2</sub>) = (1, 0), μ<sub>A</sub>(x<sub>3</sub>) = (0.5, 0.5), μ<sub>A</sub>(x<sub>4</sub>) = (0.5, 0.5), μ<sub>A</sub>(x<sub>5</sub>) = (0, 1), μ<sub>A</sub>(x<sub>6</sub>) = (0, 1)。

然后, 根据分布差别矩阵的定义, 给出表 1 所对应的差别矩阵, 如表 2 所列。按属性重要度排序后的分布差别矩阵如表 3 所列。

表 2 表 1 所对应的分布差别矩阵

|                | 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>1</sub> | ∅              | ∅              | {b, c}         | {b, c}         | {b, c, d}      | {a, b, d}      |
| x <sub>2</sub> |                | ∅              | {b, c, d}      | {b, c, d}      | {b, c}         | {a, b, d}      |
| x <sub>3</sub> |                |                | ∅              | ∅              | {d}            | {a, b, c, d}   |
| x <sub>4</sub> |                |                |                | ∅              | {d}            | {a, b, c, d}   |
| x <sub>5</sub> |                |                |                |                | ∅              | ∅              |
| x <sub>6</sub> |                |                |                |                |                | ∅              |

表 3 按属性重要度排序后的分布差别矩阵

|                | 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>1</sub> | ∅              | ∅              | {b, c}         | {b, c}         | {b, c, d}      | {b, d, a}      |
| x <sub>2</sub> |                | ∅              | {b, c, d}      | {b, c, d}      | {b, c}         | {b, d, a}      |
| x <sub>3</sub> |                |                | ∅              | ∅              | {d}            | {b, c, d, a}   |
| x <sub>4</sub> |                |                |                | ∅              | {d}            | {b, c, d, a}   |
| x <sub>5</sub> |                |                |                |                | ∅              | ∅              |
| x <sub>6</sub> |                |                |                |                |                | ∅              |

最后, 根据传统差别信息树算法给出差别信息树(见

图1),并通过算法1给出基于核属性剪枝策略的改进差别信息树(见图2),以及基于核属性剪枝和属性重要度策略的改进差别信息树(见图3),并相互比较。

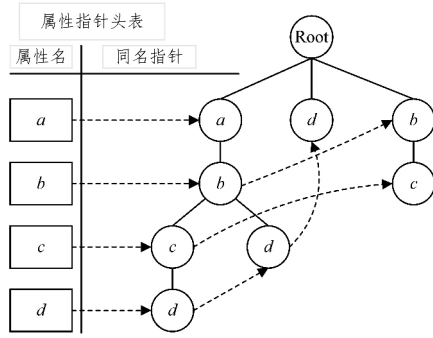


图 1

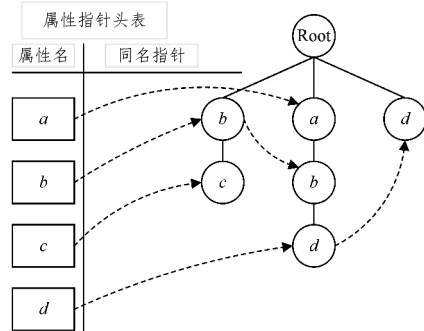


图 2

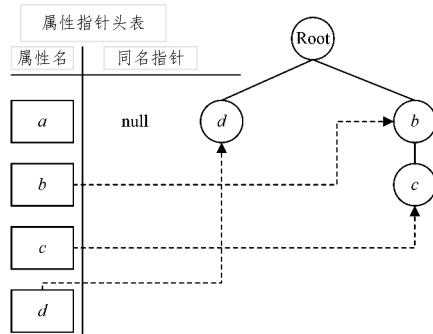


图 3

通过比较可以发现:未改进前的差别信息树(见图1)的根节点有3条子树,总共有9个节点。其中,节点b又分为2条子树,结构比较复杂而且占用空间大,冗余的信息多。因此,我们对差别信息树进行改进,提出了基于核属性剪枝策略的差别信息树(见图2)。该方法的核心思想是找到核属性元素,当构建出核属性元素的单节点后,删除之后包含核属性元素的差别信息。这个案例分析中,只用到了 $\{b,c\}, \{b,c\}, \{b,c,d\}, \{a,b,d\}, \{b,c,d\}, \{b,c,d\}, \{b,c\}, \{a,b,d\}, \{d\}$ 这9条差别信息参与差别信息树的构建,因为 $\langle d \rangle$ 已经构建出核属性的元素的单节点,之后包含元素 $d$ 的差别信息将不参与构建。图2中总共只有7个节点,根节点的子树只有3条。它完全消除了核属性的父集差别信息,而且不需要所有差别信息参与。这使差别信息树更加简洁,所需要的运算也更少。但是核属性剪枝策略只能消除核属性的父集差别信息,不能确保其他父集差别信息也能得到消除。于是,在核属性剪枝策略的基础上,我们采用属性重要度策略(见图3),以实现改

进差别信息树的最大压缩能力。在参与构建差别信息树的信息 $\{b,c\}, \{b,c\}, \{b,c,d\}, \{a,b,d\}, \{b,c,d\}, \{b,c,d\}, \{b,c\}, \{a,b,d\}, \{d\}$ 中,计算各个属性的重要度并进行排序,结果为 $\{b,c,d,a\}$ 。以属性重要度为启发信息就能加大改进差别信息树的压缩程度,并有效地消除父集差别信息。所以,图3只有2条子树,4个节点,而且不包含任何父集差别信息。

基于图2,算法2(分布属性约简算法)的求解过程如下:

- (1)初始时 $R$ 为空,而 $CList$ 为 $\{a,b,c,d\}$ ;
- (2)属性 $d$ 加入约简 $R$ 中,并从 $CList$ 删除属性 $d$ ,同时从改进差别信息树中删除含属性 $d$ 的路径,如图3所示;
- (3)选取 $CList$ 中最右边属性 $c$ ,通过同名指针链删除图1中包含属性 $c$ 的叶子节点,如图4所示;
- (4)属性 $b$ 加入约简 $R$ 中,并从 $CList$ 删除属性 $b$ ,同时从改进差别信息树中删除含属性 $b$ 的路径;
- (5)此时改进差别信息树为只含根节点的空树,算法结束,输出约简 $R = \{b,d\}$ 。

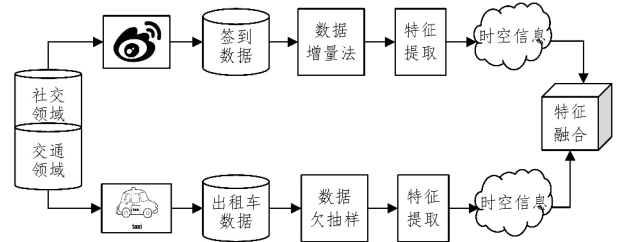


图 4

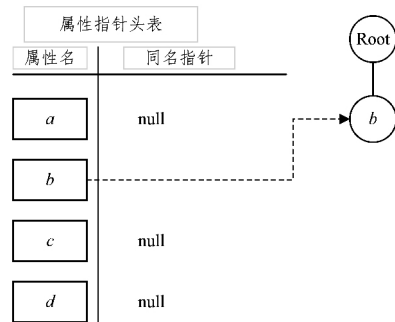


图 5

本案例分析中差别信息树的构建步骤如下:

- (1)创建差别信息树的根节点,对于决策表1中所有对象,计算差别信息,得出差别矩阵(见表2)。
  - (2)构建第一条差别信息 $\{b,c\}$ 所对应的路径 $[b,c]$ 。
  - (3)由于第二条差别信息 $\{b,c\}$ 与差别信息树中 $[b,c]$ 所对应的差别信息相同,因此第二条差别信息 $\{b,c\}$ 也映射到 $[b,c]$ 上。
  - (4)第三条差别信息 $\{b,c,d\}$ ,由于该差别信息完全包含差别信息树中路径 $[b,c]$ 对应的差别信息,采用不扩展路径策略,不构建新节点 $\langle d \rangle$ 。
  - (5)第四条差别信息 $\{a,b,d\}$ 构建第二条差别信息 $[a,b,d]$ 。
  - (6)同理,重复上面构建过程,直至把最后一条差别信息 $\{a,b,c,d\}$ 插入差别信息树中。
- 本案例分析中基于核属性剪枝策略的改进差别信息树的构建步骤如下:
- (1)重复差别信息树的构建过程中的步骤(1)一步骤(5)。
  - (2)处理第九条差别信息 $\{d\}$ 时,创建单节点 $\langle d \rangle$ ,因为

{ $d$ }是核属性元素,删除之后包含元素  $d$  的差别信息。

(3)删除之后,因为差别信息都为空集,所以构建结束。

本案例分析中基于属性重要度的改进差别信息树的构建步骤如下:

(1)基于核属性剪枝策略,得出参与构建差别信息树的所有差别信息,然后计算属性重要性程度,得出属性重要度降序排列( $b, c, d, a$ )。

(2)以属性重要度对所有差别信息进行排序,得到新的差别矩阵(见表 3)。

(3)用新的差别矩阵中的差别信息重复基于核属性剪枝策略的改进差别信息树的构建过程。

(4)直至把最后一条差别信息插入差别信息树中。

结束语 本文在差别信息树的基础上,引入核属性和属性重要度,提出了改进差别信息树。该树有着更好的压缩储存能力。但是,本文没有讨论不同的属性重要度顺序对压缩储存差别信息的作用。因此,还需要继续对不同的属性重要度进行讨论,得出一种快速属性重要度算法,并研究在不同属性重要度的情况下,进一步对差别信息压缩储存的方法。

### 参 考 文 献

- [1] 卢鹏,肖健梅,王锡淮,粗糙集属性约简的图论方法[J]. 计算机科学,2012,39(2):250-254.
- [2] 孙兴波,杨平先,干树川. 基于属性重要度的启发式特征选取算法[J]. 自动化与仪器仪表,2005(5):13-14,17.
- [3] 徐伟华. 序信息系统与粗糙集[M]. 北京:科学出版社,2013:28-32.
- [4] 李京政,杨习贝,窦慧莉,等. 重要度集成的属性约简方法研究[J]. 智能系统学报,2018,5(9):1-8.
- [5] MENG Z, SHI Z. On quick attribute reduction in decision-theoretic rough set models[J]. Information Sciences, 2016, 330(C): 226-244.
- [6] 陈昊,杨俊安,庄镇泉. 变精度粗糙集的属性核和最小属性约简算法[J]. 计算机学报,2012,35(5):1011-1017.
- [7] 王国胤,姚一豫,于洪. 粗糙集理论与应用研究综述[J]. 计算机学报,2009,32(7):1229-1246.
- [8] 张晓燕,徐伟华,张文修. 序目标信息系统中分布约简的矩阵算法[J]. 理工大学学报,2010,24(3):56-61.
- [9] 蒋云良,杨章显,刘勇. 不协调信息系统快速属性分布约简方法[J]. 自动化学报,2012,38(3):382-388.
- [10] PANG J, ZHANG X, XU W. Attribute Reduction in Intuitionistic Fuzzy Concept Lattices[J]. Abstract and Applied Analysis, 2013(9):1-13.
- [11] 于海燕,乔晓东. 一种完备的最小属性约简方法[J]. 计算机工程,2012,38(4):46-48.
- [12] 黄治国,王加阳,罗安. 一种基于分布约简的规则获取方法[J]. 计算机应用研究,2007,24(6):42-44.
- [13] XU W, LI W, LUO S. Knowledge reductions in generalized approximation space over two universes based on evidence theory[J]. Journal of Intelligent & Fuzzy Systems, 2015, 28(6): 2471-2480.
- [14] 汪凌. 不协调决策信息系统的知识约简及决策规则优化研究[J]. 计算机应用研究,2019(7):1-6.
- [15] 蒋瑜. 基于差别信息树 Rough Set 属性约简算法[J]. 控制与决策,2015,30(8):1531-1536.
- [16] JU H, YANG X, YANG P, et al. A Moderate Attribute Reduction Approach in Decision-Theoretic Rough Set [M]. Rough Sets, Fuzzy Sets, Data Mining, and Granular Computing. Springer International Publishing, 2015.
- [17] XU W, LI Y, LIAO X. Approaches to attribute reductions based on rough set and matrix computation in inconsistent ordered information systems[J]. Knowledge-Based Systems, 2012, 27(3): 78-91.
- [18] YING, HE, DAN, et al. Discernibility Matrix-Based Attribute Reduction Algorithm of Decision Table[J]. Advanced Materials Research, 2012, 457-458: 1230-1234.
- [19] 尹继亮,张楠,童向荣,等. 不协调区间值决策系统的最大分布约简[J]. 智能系统学报,2018,5(9):1-11.
- [3] PASSOS A, KUMAR V, MCCALLUM A. Lexicon Infused Phrase Embeddings for Named Entity Resolution[C]//Proceeding of the Eighteenth Conference on Computational Language Learning, 2014: 78-86.
- [4] CHIU J P C, NICHOLS E. Named Entity Recognition with Bidirectional LSTM-CNNs[J]. ArXiv:1511.08308.
- [5] COLLOBERT R, WESTON J, KARLEN M, et al. Natural Language Processing(Almost) from Scratch[J]. Journal of Machine Learning Research, 2011, 12(1): 2493-2537.
- [6] 冯艳红,于红,孙庚,等. 基于 BLSTM 的命名实体识别方法[J]. 计算机科学,2018,45(2):261-268.
- [7] 王蕾. 基于神经网络的中文命名实体识别研究[D]. 南京:南京师范大学,2017.
- [8] MNH V, HEES N, GRAVES A, et al. Recurrent models of visual attention[C]//Proceedings of the 27th International Conference on Neural Information Processing System. 2014: 2204-2212.
- [9] LUONG M T, PHAM H, MANNING C D. Effective Approaches to Attention-based Neural Machine Translation[J]. ArXiv: 1508.04025.
- [10] VASWANI A, SHAZEER N, PARMAR N, et al. Attention Is All You Need[J]. arXiv:1706.03762.
- [11] TAN Z, WANG M, XIE J, et al. Deep Semantic Role Labeling with Self-Attention[J]. ArXiv:1712.01586.
- [12] 谢志宁. 中文命名实体识别算法研究[D]. 杭州:浙江大学,2017.
- [13] GUL K S Q, 尹继泽,潘丽敏,等. 基于深度神经网络的命名实体识别方法研究[J]. 信息安全学报,2017(10):29-35.

(上接第 114 页)