

概率差别矩阵与不完备信息系统属性约简^{*})

闫德勤

(辽宁师范大学计算机系 大连 116029)

摘要 差别矩阵的概念是基于粗糙集理论对信息系统进行属性约简的一个重要内容。针对不完备信息系统的属性约简本文提出了一种概率差别矩阵的概念与构造方法,给出了相关的定理。在此基础上提出了一种利用概率差别矩阵对不完备信息系统属性约简的方法,并给出了应用举例。

关键词 粗糙集, 概率差别矩阵, 属性约简

Probability Discernibility Matrix and Attribute Reduction for Incomplete Information Systems

YAN De-Qin

(Department of Computer Science, Liaoning Normal University, Dalian 116029)

Abstract In this paper, the attribute reduction for incomplete information systems is approached. A new concept of probability discernibility matrix is proposed, and the construction method of probability discernibility matrix is given. By use of probability discernibility matrix, a new method of attribute reduction for incomplete information systems is proposed, and, as an application of this method an example is given.

Keywords Rough sets, Probability discernibility matrix, Attribute reduction

1 引言

由于不完备信息系统中一些属性的丢失或不确定,使得以等价类为基础的针对完备信息系统属性约简的粗糙集理论方法受到限制^[5,6]。为对不完备信息系统进行属性约简,一些学者对信息系统中丢失属性或不确定数据从不同的角度提出了一些处理方法^[5],也有一些对传统粗糙集扩展模型的提出。所有这些方法与模型都是不同意义下对不完备信息系统进行属性约简的一种处理方式。差别矩阵的概念是基于粗糙集理论对信息系统进行属性约简的一个重要内容,本文提出了一种概率差别矩阵的概念与构造方法,给出了相关的定理。在此基础上提出了一种利用概率差别矩阵对不完备信息系统属性约简的方法,并给出了应用举例。

2 基本概念

设 $S=(U, Q, V, F)$ 为一信息系统,其中 $U=\{x_1, x_2, \dots, x_n\}$ 是论域, Q 是属性集合, V 是属性取值集合, F 是 $U \times Q \rightarrow V$ 的映射。设属性集合包含 m 个条件属性 $C=\{C_1, C_2, \dots, C_m\}$ 和一个决策属性 D 。若 D 的取值有 s 个,则由 D 导出的等价类构成 U 的一个划分: $\{Y_1, Y_2, \dots, Y_s\}$ 。其中, $Y_i=\{x \in U | F(x, D)=i\}, i=1, 2, \dots, s$ 。

信息系统中每一 x_i 及其所对应的属性值称为一个(信息表示的)规则。本文中简称 x_i 规则。

在一个信息系统中,当 $i \neq j$ 时若存在 $F(x_i, C)=F(x_j, C)$ 但 $F(x_i, D) \neq F(x_j, D)$ 则称该系统为不相容的,此时 x_i 与 x_j 所对应的规则为不相容规则。

定义 1 设 $X \subseteq U$ 为论域的一个子集, $P \subseteq C$, X 的关于 P 的下近似为

$$P_-X = \{x \in U | [x]_P \subseteq X\}$$

其中 $[x]_P$ 表示 U 中在等价关系 P 下的等价类元素构成的集合。

定义 2 设 U 为一个论域, P, Q 为 U 上的两个等价关系

簇, Q 的 P 正域记为 $POS_P(Q)$, 定义为

$$POS_P(Q) = \bigcup_{x \in U; Q} P_-(X)$$

定义 3 设 $P \subseteq C$, 对于划分 $\{Y_1, Y_2, \dots, Y_k\}$ 的 P 的近似精度为

$$\gamma_P = \frac{\sum_{i=1}^k \text{card}(P_-Y_i)}{\text{card}(U)}$$

其中, $\text{card}()$ 表示集合的基数。

定义 4 设 $P \subseteq C$, 若 $\gamma_P = \gamma_C$, 且不存在 $R \subset P$, 使得 $\gamma_R = \gamma_P$, 则称为 P 为 C 的一个属性约简。所有 C 的属性约简的交称为 C 的核, 记为 $Core(C)$ 。

定义 5^[4] 给定信息系统 S , 差别矩阵 $M=(m_{ij})$ 的元素定义为: 当 $F(x_i, D) \neq F(x_j, D)$ 时, $m_{ij} = \{a \in C | F(x_i, a) \neq F(x_j, a)\}$, 在其它情况下 m_{ij} 为空集。

文[4]指出当差别矩阵中的某个元素为单元素时, 该属性属于核。

若信息系统中的每个属性值都是已知的, 则称为完备的信息系统。在有些情况下, 由于种种原因信息系统中的某些属性值不能得到或确定, 则称信息系统为不完备信息系统。

3 概率差别矩阵及相关定理

对于完备信息系统, 可以利用差别矩阵进行属性约简^[4,6], 而对于不完备信息系统由于未知属性值的存在不能采用已有的方法构造差别矩阵^[5]。为利用差别矩阵对不完备信息系统进行属性约简, 下面提出一种关于不完备信息系统的差别矩阵构造方法, 以该方法构造的差别矩阵称为概率差别矩阵。

设 $S=(U, Q, V, F)$ 为一个不完备信息系统, V 是属性取值集合, 信息系统中的不确定属性值用 $*$ 表示, 系统中属性取值的种类用 v 表示。构造概率差别矩阵基于不确定属性值取到信息系统中任意一个确定属性值的可能性为 v 分之一的概率假设。

定义 6 给定不完备信息系统 S , 概率差别矩阵 $M=$

^{*}) 国家自然科学基金(60372071)资助; 辽宁师范大学校基金资助。闫德勤 博士, 教授, 主要研究领域为模式识别、数据挖掘等。

一种新型的神经网络构造方法 RCBNN^{*})

谢振华 李宁 商琳 陈兆乾 陈世福

(南京大学计算机软件新技术国家重点实验室 南京 210093)

摘要 一个好的神经网络结构可以大大提高它的处理能力和收敛速度,所以神经网络的构造方法一直是人们研究的热点问题。本文利用粗集理论的数据分析能力和决策树对数值属性的分割能力,提出一种基于粗集与决策树的新型神经网络构造方法 RCBNN。经试验表明,使用该方法构造的神经网络,具有易于构造、可理解性好、收敛速度快且构造的网络规模较小的特点。

关键词 神经网络,粗集,决策树

RCBNN: A New Constructing Method for Neural Network

XIE Zhen-Hua LI Ning SHANG Ling CHEN Zhao-Qian CHEN Shi-Fu

(State Key Laboratory for Novel Software Technology, Nanjing University, Nanjing 210093)

Abstract A good structure can greatly improve the power and convergent speed of neural network. According this, the constructing method is always a hotspot in neural network research. This paper proposes a constructing method for neural network based on Rough Set and Decision Tree. Rough Set has good capability of data analysis; Decision Tree is good at segmentation of continues-valued attributes. This method makes good use of these advantages. Experimental results show that the neural network designed by this method is easy constructible, good understandable, fast convergent, and small dimensional.

Keywords Neural network, Rough set, Decision tree

1 引言

波兰科学家 Pawlak 提出的粗集(Rough set)^[1]是描述不完整、不精确和含噪声数据的有力工具,其主要特点是不需要任何与数据有关的先验和附加知识。从本质上看,它反映了认知过程在非确定性、非模型化信息处理方面的机制和特点,从而成为一种有效的非单调推理工具。粗糙集模型是一种结构化的、非数值化的信息处理方法,处理的是数据的符号描述,以属性、语义决策规则等形式构造知识表达,其主要特点是不需要任何与数据有关的先验和附加知识,这一点明显区别于模糊集理论中的隶属函数等方法。目前,粗糙集主要应用于属性约简、规则生成及预测等几个方面。但粗集理论容错能力与推广能力相对较弱,适于处理离散的非确定数据,对于连续数据的处理能力有限^[2],且因对数值型属性作为离散看待,规则生成后的规模大;决策树学习是应用最广的归纳推理算法之一,对噪声数据有很好的健壮性,其学习速度快,规则学习能力强,本文所采用的 C4.5 决策树算法^[3]能通过

对连续数据的分割,有效地缩小规则的规模,但决策树容易产生过拟合现象,虽然通过裁减可一定程度避免过拟合,但精度的损失较大;人工神经网络最具吸引力的是其数值逼近能力,并能够处理定量的、数值化的信息。人工神经网络关注的是建模和学习过程中获得的数据的细节信息,具有较强的自组织能力、容错能力和推广能力,但不能优选条件属性,而粗集的属性约简能力恰好能弥补这一缺陷。已有研究者成功地将粗集理论与神经网络结合,构造出新型的神经网络结构^[4,5]。本文利用粗集与决策树理论的优点,并结合人工神经网络的特点提出了一种基于粗集与决策树的神经网络构造方法——RCBNN,该方法易于构造、可理解性好、计算简单、收敛速度快且构造的神经网络规模较小,通过实验表明,该方法能够得到较满意的结果。

2 粗集理论与决策树算法

给定一个有穷对象集 U 、有穷属性集 A 、各属性值的集合 V 以及信息函数 f 构成一个信息系统,表示为 (U, A, V, f) ,

*)基金项目:国家自然科学基金(60273033);江苏省自然科学基金(BK2003067)。谢振华 硕士研究生,主要研究方向:神经计算。李宁 硕士,讲师,主要研究方向:机器学习。商琳 博士研究生,主要研究方向:人工智能、数据挖掘。陈兆乾 教授,博士生导师,主要研究领域:人工智能、机器学习。陈世福 教授,博士生导师,主要研究领域:人工智能、机器学习。

及 Vague 集属性信息系统的属性约简,其结果将另文发表。

参 考 文 献

- 1 Pawlak Z. Rough set approach to multi-attribute decision analysis. *European Journal of Operational Research*, 1994, 72: 443~459
- 2 Pawlak Z. Rough set theory and its applications to data analysis. *Cybernetics and System*, 1998, 29(27): 661~688
- 3 Chen M S, Han J, Yu P S. Data Mining: An overview from a database perspective. *IEEE Transaction on Knowledge and Data Engineering*, 1996, 8(6): 866~883
- 4 Hu XiaoHua, Cercone N. Learning in relational database: a rough set approach. *Computational Intelligence*, 1995, 11(2): 323~337
- 5 Grzymala-Busse J W, Hu M. A comparison of several approaches to missing attribute values in data mining. In: *Proc. of the 2nd Int'l. Conf. on Rough Sets and Current Trends in Computing*. Berlin: Springer Verlag, 2000. 378~385
- 6 王国胤. 粗糙集理论与知识获取. 西安:西安交通大学出版社, 2001
- 7 曾黄麟. 粗糙集理论及其应用. 重庆:重庆大学出版社, 1996