

一种基于决策粗糙集的模糊C均值聚类数的确定方法

石文峰 商琳

(南京大学计算机科学与技术系 南京 210023)

(计算机软件新技术国家重点实验室(南京大学) 南京 210023)

摘要 Fuzzy C-Means(FCM)是模糊聚类中聚类效果较好且应用较为广泛的聚类算法,但是其对初始聚类数的敏感性导致如何选择一个好的C值变得十分重要。因此,确定FCM的聚类数是使用FCM进行聚类分析时的一个至关重要的步骤。通过扩展决策粗糙集模型进行聚类的有效性分析,并进一步确定FCM的聚类数,从而避免了使用FCM时不好的初始化所带来的影响。文中提出了一种基于扩展粗糙集模型的模糊C均值聚类数的确定方法,并通过图像分割实验来验证聚类的效果。实验通过比对不同聚类数下分类结果的代价获得了一个较好的分割结果,并将结果与Z. Yu等人于2015年提出的蚁群模糊C均值混合算法(AFHA)以及提高的AFHA算法(IAFHA)进行对比,结果表明所提方法的聚类结果较好,图像分割效果较明显,Bezdek分割系数比AFHA和IAFHA算法的更高,且在Xie-Beni系数上也有较大优势。

关键词 模糊C均值,决策粗糙集,图像分割

中图分类号 TP181 文献标识码 A DOI 10.11896/j.issn.1002-137X.2017.09.008

Determining Clustering Number of FCM Algorithm Based on DTRS

SHI Wen-feng SHANG Lin

(Department of Computer Science & Technology, Nanjing University, Nanjing 210023, China)

(State Key Laboratory for Novel Software Technology at Nanjing University, Nanjing 210023, China)

Abstract Fuzzy C-Means(FCM), as the most popular algorithm of the soft clustering, has been extensively used to make compact and well separated clusters. However, its sensitivity to initial cluster number makes choosing a better C value become very important. So it is an important step to determine the number of FCM clustering when we use FCM to do cluster analysis. In this paper, the extended decision-theoretic rough sets(DTRS) model is applied for the purpose of clustering validity analysis which could overcome the defect of the FCM algorithm. We proposed the method for determining clustering number of FCM algorithm based on DTRS, and we verified the effect of the clustering by image segmentation. Good segmentation results can be obtained when we compare the cost of different number of clusters. We compared our results with the ant colony fuzzy c-means hybrid algorithm (AFHA), which was proposed by Z. Yu et al in 2015, and the improved AFHA (IAFHA). The experimental results show that our clustering result is better in Bezdek partition coefficient with a higher value than AFHA and IAFHA algorithms, and in the Xie-Beni index as well.

Keywords Fuzzy C-Means, Decision-theoretic rough sets, Image segmentation

聚类已经被广泛应用到诸如模式识别、图像处理以及医疗诊断等许多领域。聚类策略可以大体分成两大类:硬聚类划分和软聚类划分。硬聚类需要每个对象对于某一个类别的隶属度非0即1;与其不同,软聚类将隶属度的概念加入聚类方法中,使得每个对象属于某个特定类的隶属度可以在0到1之间取值,这通常和实际应用是一致的,因为众多情况下对象之间并没有明确的界限。

在软聚类分析中,最受大众欢迎的即是模糊C均值FCM

(Fuzzy C-Means)算法,其最早由Dunn^[1]于1973年提出并由Bezdek^[2]于1981年做了进一步的改进。该算法基于最小化的目标函数原理,通过更新隶属度函数和聚类中心来不断迭代产生最后的结果。尽管FCM算法比硬聚类算法的效果好,但是它仍然存在聚类中心和聚类数需要被提前确定的缺陷。一个好的初始化能够使得最终的聚类效果较为理想,而一个不合适的初始化可能会导致较差的聚类效果。鉴于FCM对聚类数初始化的严重依赖,本文主要是讨论如何能够

更好地初始化 FCM,即确定一个较为合适的聚类数。

为了得到较为合适的聚类数,使用扩展的决策粗糙集(DTRS)模型。传统的粗糙集理论是由 Pawlak^[3]于 1982 年提出的,它采用上近似集合和下近似集合来描述一个粗糙集合。然而,Pawlak 粗糙集模型并未考虑决策规则的错误容忍度。鉴于这一不足,Yao 等人^[4]于 20 世纪 90 年代提出了一种 DTRS 模型,它是对经典 Pawlak 粗糙集理论的概率拓展,并且提供了一个对分类的更好理解。该模型通过引入贝叶斯理论来最小化损失代价。最近,DTRS 模型被广泛应用到众多领域,如信息过滤^[5]、属性约简^[6]、投资决策^[7]、垃圾邮件过滤^[8]等。在聚类方面,Lingras 等人^[9]描述了如何在 DTRS 模型中通过调整损失函数来提高聚类的有效性指标。Yu 等人^[10]则通过扩展 DTRS 模型提出了一个分层方法来进行聚类有效性的分析。为了能够获得更好的聚类结果,提出一种基于决策粗糙集的模糊 C 均值聚类数的确定方法,其通过使用 DTRS 进行聚类有效性分析,辅助 FCM 得到较为合适的聚类结果。

近年来,很多通过聚类有效性分析来克服 FCM 较依赖初始化的方法已经被提出。文献[12]提出了一种蚁群模糊 C 均值混合算法(AFHA),这一算法使用蚁群算法(AS)完成 FCM 的初始化。为了提高 AFHA 的计算效率,文献[12]又提出了一个提高的 AFHA(IAFHA)算法。然而,文献[12]中的算法存在参数过于冗余的不足。

本文提出了一种通过 DTRS 来辅助 FCM 确定聚类数从而获得较好聚类效果的方法——DTRS&FCM。第 1 节将简述模糊 C 均值和决策粗糙集的相关原理;第 2 节描述基于模糊 C 均值和决策粗糙集的算法的原理和流程;第 3 节给出了实验结果及其分析;最后总结全文,并指出进一步的工作方向。

1 模糊 C 均值和扩展的决策粗糙集

1.1 模糊 C 均值(FCM)

模糊聚类是一种每个对象可以部分属于某个类,即隶属度取值不仅限于 0 和 1 两个值,而是在 $[0,1]$ 区间内取值的一种聚类算法。也就是说,对象对于类不再具有非此即彼的特性,而可以同时属于多个类,具体隶属于类的程度可以用隶属度函数来表示。在多种模糊聚类算法中,最重要也最广泛使用的是模糊 C-均值(Fuzzy C-Means)算法。该算法最早是由 Dunn^[1]将硬聚类的聚类准则函数推广到模糊聚类中而提出的带有隶属度加权的误差平方和函数,进而由 Bezdek^[2]做了进一步的扩展工作。

标准的 FCM 算法是一种聚类方法中的迭代方法,它通过最小化加权的误差平方和将一组对象聚类成最佳的 C 类。该算法使得各个类之间可以有重叠的部分,允许对象被包含在多个类中。该算法需要最小化如下目标函数:

$$J_m = \sum_{k=1}^c \sum_{i=1}^n u_{ki}^m \|x_i - v_k\|^2 \quad (1)$$

$$\text{s. t. } \sum_{k=1}^c u_{ki} = 1, u_{ki} \in [0, 1], 0 \leq \sum_{i=1}^n u_{ki} \leq n \quad (2)$$

其中, n 是数据对象的个数; $c(2 \leq c \leq n)$ 是聚类的个数;运算符 $\|$ 为欧氏范数; x_i 是 d 维向量空间的第 i 个数据对象;参数 $m(m \geq 1)$ 是 FCM 的隶属度函数权重因子,它可以控制划分数据的模糊程度, m 的值越大则函数的模糊性越大; $v_k(1 \leq k \leq c)$ 表示类中心; $u_{ki}(1 \leq k \leq c, 1 \leq i \leq n)$ 表示对象 x_i 隶属于第 k 个类的程度; $U = \{u_{ki}\}$ 是 $c \times n$ 的矩阵; $V = \{v_1, v_2, \dots, v_c\}$ 是 $s \times c$ 的矩阵。

FCM 的算法框架如算法 1 所示。

算法 1 模糊 C 均值(FCM)

Step1 初始化:设定迭代终止阈值 ϵ ,聚类数 C ,模糊因子 m ;

Step2 设置迭代计数器 $q=0$;

Step3 随机初始化类中心矩阵:

$$C^{(q)} = c_k (k=1, 2, \dots, c)$$

Step4 根据 $C^{(q)}$ 计算隶属度矩阵 $U^{(q)}$:

$$u_{ji} = \frac{1}{\sum_{k=1}^c \left(\frac{d_{ji}}{d_{ki}}\right)^{2/m-1}} \quad (3)$$

注意,如果 $d_{ji}=0$,则 $u_{ji}=1, u_{li}=0(1=1, \dots, N \text{ and } l \neq j)$

Step5 根据 $U^{(q)}$ 计算类中心矩阵 $C^{(q+1)}$:

$$c_j = \frac{\sum_{i=1}^n u_{ji}^m x_i}{\sum_{i=1}^n u_{ji}^m} \quad (4)$$

Step6 根据 $C^{(q+1)}$ 以及式(3)更新隶属度矩阵 $U^{(q+1)}$;

Step7 如果 $\max\{U^{(q+1)} - U^{(q)}\} \leq \epsilon$,则算法终止;否则, $q=q+1$,返回 Step4。

FCM 是在传统 K-Means 算法上的延伸,它通过最小化类内间距不断更新隶属度。当所有对象对所有类的隶属度都确定完毕后,对于每个对象,选择隶属度最大的类作为该对象所属的类别。

1.2 扩展的决策粗糙集(DTRS)

基于 Yao 等人^[4]于 20 世纪 90 年代提出的一种 DTRS 模型,Yu 等人^[5]提出了基于扩展决策粗糙集模型在层次聚类算法框架下自动确定的聚类数的方法。该理论不再对单独的对象而言,而是对对象对 (x_i, x_j) 的分类状态而言,扩展了决策粗糙集模型中的相关定义。

设 $\Omega = \{C, \neg C\}$ 为一个两互补的状态集合,其中 C 和 $\neg C$ 分别表示两个对象 x_i 和 x_j 属于同一个类别的状态和两个不同类别的状态; $A = \{a_P, a_N\}$ 表示一个包含两种可能的行动集合, a_P 和 a_N 分别表示两个对象 x_i 和 x_j 被分到相同类别的决策; $P(C|(x_i, x_j))$ 表示将对象 x_i 和 x_j 分配到状态 C 的概率; $sim(x_i, x_j)$ 表示对象 x_i 和 x_j 之间的相似度; $\lambda_{a_j w_j}$ 表示处于状态 w_j 时做出决策 a_j 后的代价损失。

显而易见, $P(C|(x_i, x_j))$ 与对象 x_i 与 x_j 之间的相似度成正比。为了建立 $P(C|(x_i, x_j))$ 和 $sim(x_i, x_j)$ 之间的关系,引进阈值 val ,如果 $sim(x_i, x_j) = val$,则 $P(C|(x_i, x_j)) = 0.5$ 。这里使用如下公式计算 val :

$$val = \frac{1}{N^2} \cdot \sum_{i=1}^N \sum_{j=1}^N sim(x_i, x_j) \quad (5)$$

从而 $P(C|(x_i, x_j))$ 和 $P(\neg C|(x_i, x_j))$ 按如下公式计算:

$$P(C|(x_i, x_j)) = \begin{cases} 0.5 + \frac{sim(x_i, x_j) - val}{2 - 2val}, & sim(x_i, x_j) \geq val \\ 0.5 - \frac{val - sim(x_i, x_j)}{2val}, & sim(x_i, x_j) < val \end{cases} \quad (6)$$

$$P(\neg C|(x_i, x_j)) = 1 - P(C|(x_i, x_j)) \quad (7)$$

接下来,用如下公式计算将两个对象 x_i 和 x_j 分到同一类和不同类的期望代价损失:

$$R(a_P|(x_i, x_j)) = \lambda_{PP}P(C|(x_i, x_j)) + \lambda_{PN}P(\neg C|(x_i, x_j)) \quad (8)$$

$$R(a_N|(x_i, x_j)) = \lambda_{NP}P(C|(x_i, x_j)) + \lambda_{NN}P(\neg C|(x_i, x_j)) \quad (9)$$

其中,我们只考虑了一种代价函数,即两个属于同一类的对象被划分到同一类中的代价为 0;两个不属于同一类的对象被划分到同一类中的代价为 1;反之类似。用损失函数可以表示为:

$$\begin{aligned} \lambda_{PP} &= 0, \lambda_{PN} = 1 \\ \lambda_{NP} &= 1, \lambda_{NN} = 0 \end{aligned}$$

则式(8)可以改写为:

$$\begin{aligned} R(a_N|(x_i, x_j)) &= \lambda_{NP}P(C|(x_i, x_j)) \\ R(a_P|(x_i, x_j)) &= \lambda_{PN}P(\neg C|(x_i, x_j)) \end{aligned} \quad (10)$$

所以,对于一组对象 x_i 和 x_j 而言,在聚类结果 CS_i 的状态下的代价损失为:

$$R(CS_i|(x_i, x_j)) = \begin{cases} P(C|(x_i, x_j)), & (x_i, x_j) \in C \\ P(\neg C|(x_i, x_j)), & (x_i, x_j) \in \neg C \end{cases} \quad (11)$$

将式(6)和式(7)代入到式(11)中,有:

$$R(CS_i|(x_i, x_j)) = \begin{cases} 0.5 + \frac{sim(x_i, x_j) - val}{2 - 2val}, & sim(x_i, x_j) > val, (x_i, x_j) \in C \\ 0.5 - \frac{val - sim(x_i, x_j)}{2val}, & sim(x_i, x_j) \leq val, (x_i, x_j) \in C \\ 0.5 + \frac{sim(x_i, x_j) - val}{2 - 2val}, & sim(x_i, x_j) > val, (x_i, x_j) \in \neg C \\ 0.5 - \frac{val - sim(x_i, x_j)}{2val}, & sim(x_i, x_j) \leq val, (x_i, x_j) \in \neg C \end{cases} \quad (12)$$

为了评估聚类结果的好坏,定义评估函数为:

$$R(CS_i) = \sum_{i=1}^n \sum_{j=1}^n R(CS_i|(x_i, x_j)) \quad (13)$$

2 基于决策粗糙集的模糊 C 均值聚类数的确定方法

损失函数是用来描述做出不同决策所带来的不同损失,若同一个类的数据被分到不同类或不同类的数据被划分到同一类,则代价损失较大;反之,则代价损失较小。我们使用式(13)作为判断聚类结果的损失函数。

同时,对于 $\lambda_{PP}, \lambda_{PN}, \lambda_{NP}$ 和 λ_{NN} 即由相同类划到相同类和不同类的代价以及由不同类划到相同类和不同类的代价,不采用上节所提到的 0, 1 值,而是根据实际情况设定。

$sim(x_i, x_j)$ 即对象 x_i 和 x_j 之间的相似度由先验给出。一般情况下,相似度使用欧氏距离或者余弦距离计算。本算法的实验部分采用如下定义的巴氏系数来计算相似度:

$$BC(p, q) = \frac{d}{\sum_{i=1}^d \sqrt{p_i(x)q_i(x)}} \quad (14)$$

因此,本文提出的算法即给定 FCM 一个初始的 C 值,并计算聚类结果的代价,然后不断更改 C 值,若选定的 C 值越接近实际应划分的聚类数,则损失函数的值应越小,所以只需对应取损失函数最小的 C 值即可。算法流程如算法 2 所示。

算法 2 DTRS&-FCM

输入: x_i 和 x_j 的相似度矩阵; $\lambda_{PP}, \lambda_{PN}, \lambda_{NP}$ 和 λ_{NN}

输出: 最终聚类结果

Step1 使用式(5)计算阈值 val ;

Step2 使用式(6)、式(7)得到概率矩阵 $P(C|(x_i, x_j))$;

Step3 初始化 FCM 的聚类数 $C=2$;

Step4 将 C 的值传入 FCM,并使用 FCM 算法得到相应的聚类结果;

Step5 根据给定的 $\lambda_{PP}, \lambda_{PN}, \lambda_{NP}$ 和 λ_{NN} 以及式(8)、式(9)计算 $R(CS_i|(x_i, x_j))$

Step6 根据式(13)计算整个聚类结果的代价,并将其与上一个代价进行比较,若小于前一个代价,则 $C=C+1$,转 Step4;否则,输出上一个 C 值以及在上一 C 值下的聚类结果。

3 实验结果与分析

3.1 实验环境与实验参数

在 Intel(R) Core(TM) i5-4900 CPU 3.30GHz 和 Windows 10 的环境下进行实验,编程语言为 Matlab。在实验中,使用所提方法对图像进行分割,以证明所提算法可以给出正确的聚类数,并且能够得到较为优异的效果。我们选取了 UC-Berkeley 图像分割数据集^[11]的若干张图片作为图像分割的数据,并使用它们来测试所提算法。

实验中使用了式(14)即巴氏系数来计算相似度; $\lambda_{PP}, \lambda_{PN}, \lambda_{NP}$ 和 λ_{NN} 的值则分别设置为 0, 1, 1, 0;而在 FCM 中设置迭代终止阈值 ϵ 为 0.0001,模糊因子 m 为 1。

3.2 实验评估函数

实验使用了两个基准系数来评估实验结果,一个是 Bezdek 分割系数^[13],函数定义为:

$$V_{FC} = \frac{\sum_{i=1}^N \sum_{j=1}^M u_{ji}^2}{N} \quad (15)$$

这个评估函数可以用来度量分割结果的模糊性,它的取值位于 0 和 1 之间。 V_{FC} 的值越大,表示像素的类属性程度越高,聚类分割的效果越好。

另一种基准系数是 Xie-Beni 函数^[14],其定义为:

$$V_{XB} = \frac{\sum_{i=1}^N \sum_{j=1}^M u_{ji}^2 \|x_i - c_j\|^2}{N \min_{j \neq k} \{ \|c_j - c_k\|^2 \}} \quad (16)$$

一般情况下,类内间距越小,类内距离越大,表明聚类结果越好。因此, V_{XB} 值越小,聚类结果越好。

3.3 实验结果

对实验结果给出了两种呈现:1)如图 1 所示的可视化分

割结果;2)如表1和表2所列的与其他方法的基准系数的对比结果。

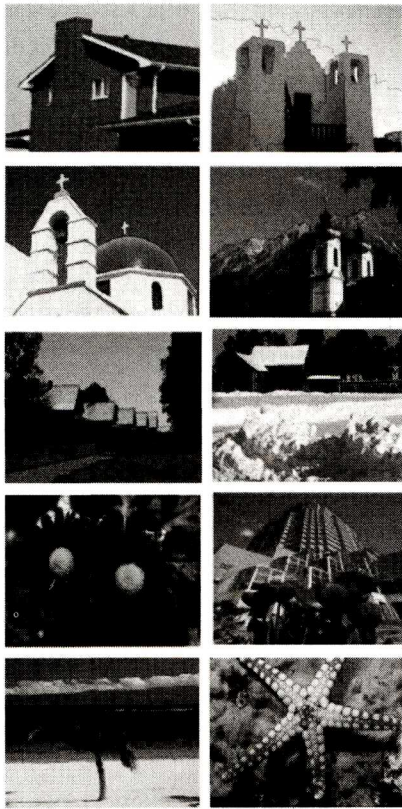


图1 DTRS&FCM算法下的图像分割效果

表1 不同算法下的 V_{FC} 结果

Images	Algorithms		
	AFHA	IAFHA	DTRS&FCM
Church3	0.586	0.586	0.794
SunFlower	0.631	0.621	0.804
Church2	0.753	0.771	0.929
Beach	0.587	0.603	0.781
Building	0.498	0.472	0.702
House	0.736	0.729	0.879
Church1	0.713	0.746	0.864
House1	0.636	0.637	0.767
House3	0.668	0.670	0.860
StarFish	0.497	0.505	0.810
Sky	0.583	0.583	0.866
Butterfly	0.478	0.480	0.812
Lena	0.468	0.453	0.737
Peppers	0.498	0.513	0.661

表2 不同算法下的 V_{XB} 结果

Images	Algorithms		
	AFHA	IAFHA	DTRS&FCM
Church3	0.214	0.272	0.178
SunFlower	0.258	0.177	0.342
Church2	0.216	0.283	0.107
Beach	0.231	0.189	0.097
Building	0.248	0.269	0.620
House	0.118	0.086	0.085
Church1	0.136	0.117	0.101
House1	0.267	0.219	0.248
House3	0.182	0.140	0.107
StarFish	0.279	0.279	0.456
Sky	0.294	0.294	0.221
Butterfly	0.437	0.350	0.320
Lena	0.316	0.284	0.479
Peppers	0.718	0.820	0.469

3.4 实验分析

从可视化结果来看,该算法的整体分割效果较好。

基于模糊聚类系数评估标准,从表1可以明显观察到,所提方法在Bezdek分割系数方面都能获得比其他方法更加理想的取值,即文中所提方法对于像素属于某个类的确定性程度会比较高,而不像其他方法具有较大的模糊性。而从表2可以观察到,Xie-Beni系数的表现力相对Bezdek系数而言不稳定性较大,但所提算法还是相对具有优势。

尽管算法效率不是我们的主要关注点,但是运行时间对算法的实用性有着较大的影响。所提算法的算法复杂度为 $O(MN^2)$ (M 为最终聚类数),而FCM和蚁群算法(AS)的算法复杂度皆为 $O(MN)$,因此所提算法的时间成本较AFHA和IAFHA稍高。

结束语 由于FCM存在着对初始数据敏感的缺陷,本文将DTRS与FCM结合,使用DTRS来计算FCM在每个聚类数下聚类结果的代价,并通过分析相应的代价成功给出了效果最佳的聚类数。如何获得FCM的最佳聚类数一直是FCM长期存在的困扰,而文中为使用FCM的研究者提供了一个可以自动确定聚类数的新方法。实验中使用本文提出的算法进行图像分割,给出了图像分割的视觉效果,并基于FCM算法中常用的Bezdek和Xie-Beni评估系数将所提算法与其他算法进行对比,从结果中可以看出所提方法是行之有效的。

参考文献

- [1] DUNN J C. A fuzzy relative of the isodata process and its use in detecting compact well-separated clusters[J]. Journal of Cybernetics, 1974, 3(3): 32-57.
- [2] BEZDEK J C. Pattern recognition with fuzzy objective function algorithms[M]. Kluwer Academic Publishers, 1981.
- [3] PAWLAK Z. Rough sets[J]. International Journal of Computer & Information Sciences, 1982, 11(5): 341-356.
- [4] YAO Y, WONG S K M. A decision theoretic framework for approximating concepts[J]. International Journal of Man-machine Studies, 1992, 37(6): 793-809.
- [5] ZHAO W Q, ZHU Y L, GAO W. Information filtering model based on decision-theoretic rough set theory[J]. Computer Engineering and Applications, 2007, 43(7): 185-187. (in Chinese)
赵文清, 朱永利, 高伟. 一个基于决策粗糙集理论的信息过滤模型[J]. 计算机工程与应用, 2007, 43(7): 185-187.
- [6] JIA X, LIAO W, TANG Z, et al. Minimum cost attribute reduction in decision-theoretic rough set models [J]. Information Sciences an International Journal, 2013, 219(1): 151-167.
- [7] LIU D, YAO Y, LI T. Three-way investment decisions with decision-theoretic rough sets[J]. International Journal of Computational Intelligence Systems, 2011, 4(1): 66-74.
- [8] JIA X, ZHENG K, LI W, et al. Three-way decisions solution to filter spam email: an empirical study[M]//Rough Sets and Current Trends in Computing. Springer Berlin Heidelberg, 2012: 287-296.

$(1, \text{高}, l_2) \rightarrow (0.5, [,])$, $(0, \text{中}, l_1) \rightarrow (0.5, [,])$
 $(0, \text{低}, l_3) \rightarrow (0.5, [,])$, $(*, \text{中}, l_1) \rightarrow (0.5, [,])$
 $(1, \text{高}, l_5) \rightarrow (0.5, [,])$, $(0, \text{中}, l_2) \rightarrow (0.5, [,])$
 $(0, \text{中}, l_4) \rightarrow (0.5, [,])$, $(*, \text{中}, l_1) \rightarrow (0.6, [,])$
 $(0, \text{中}, l_2) \rightarrow (0.6, [,])$, $(0, \text{中}, l_1) \rightarrow (0.6, [,])$
 $(1, \text{高}, l_2) \rightarrow (0.6, [,])$, $(0, \text{中}, l_4) \rightarrow (0.6, [,])$
 $(1, \text{高}, l_5) \rightarrow (0.6, [,])$, $(0, \text{中}, l_4) \rightarrow (0.7, [,])$
 $(1, \text{高}, l_5) \rightarrow (0.7, [,])$, $(0, \text{中}, l_2) \rightarrow (0.6, [6, 7])$
 $(1, \text{高}, l_2) \rightarrow (0.6, [6, 7])$, $(1, \text{高}, l_2) \rightarrow (0.6, [6, 7])$
 $(0, \text{中}, l_4) \rightarrow (0.8, [6, 8])$, $(1, \text{高}, l_4) \rightarrow (0.8, [6, 8])$
 $(1, \text{高}, l_5) \rightarrow (0.9, [6, 9])$

根据定义 10, 进一步得到如下非冗余决策规则:

$(0, \text{低}, l_1) \rightarrow (0.5, [,])$
 $(0, \text{中}, l_1) \rightarrow (0.6, [,])$
 $(0, \text{中}, l_4) \rightarrow (0.7, [,])$
 $(0, \text{中}, l_2) \rightarrow (0.6, [6, 7])$
 $(0, \text{中}, l_4) \rightarrow (0.8, [7, 8])$
 $(1, \text{高}, l_5) \rightarrow (0.9, [6, 9])$

结束语 本文就决策形式背景异构数据的情况讨论了知识发现问题, 主要给出了异构(决策)形式背景的定义, 研究了概念格构造, 给出了决策规则, 得到了非冗余规则挖掘算法。

异构数据的知识发现是一个非常重要的研究课题, 虽然本文提出了一些可行的分析方法, 但是它们的时间复杂度均是指数级的, 这意味着它们面对大规模数据时效率不高, 因此继续探讨更加有效的知识发现方法是今后的工作方向。

参考文献

- [1] WILLE R. Restructuring lattice theory: an approach based on hierarchies of concepts[M]// Rival I, ed. Ordered Sets. Dordrecht-Boston; Reidel, 1982; 445-470.
- [2] HU K Y, LU Y C, SHI C Y. Advances in concept lattice and its application[J]. Journal of Tsinghua University (Science and Technology), 2000, 40(9): 77-81. (in Chinese)
胡可云, 陆玉昌, 石纯一. 概念格及其应用进展[J]. 清华大学学报(自然科学版), 2000, 40(9): 77-81.
- [3] ZHANG W X, WEI L, QI J J. Attribute reduction theory and approach to concept lattice[J]. Science China Series F—Information Sciences, 2005, 35(6): 628-639. (in Chinese)
- [4] LI J H, MEI C L, XU W H, et al. Concept learning via granular computing: A cognitive viewpoint[J]. Information Sciences, 2015, 298: 447-467.
- [5] XU W H, LI W T. Granular computing approach to two-way learning based on formal concept analysis in fuzzy datasets[J]. IEEE Transactions on Cybernetics, 2016, 46(2): 366-379.
- [6] ZHANG T, REN H L, HONG W X, et al. The visualizing calculation of formal concept that based on the attribute topologies[J]. Acta Electronica Sinica, 2014, 42(5): 925-932. (in Chinese)
张涛, 任宏雷, 洪文学, 等. 基于属性拓扑的可视化形式概念计算[J]. 电子学报, 2014, 42(5): 925-932.
- [7] QI J J, WEI L, YAO Y Y. Three-way formal concept analysis[M]// Rough Sets and Knowledge Technology. 2014: 732-741.
- [8] WEI L, QI J J, ZHANG W X. Attribute reduction theory of concept lattice based on decision formal contexts[J]. Science China Series F—Information Sciences, 2008, 38(2): 195-208. (in Chinese)
魏玲, 祁建军, 张文修. 决策形式背景的概念格属性约简[J]. 中国科学(F辑): 信息科学, 2008, 38(2): 195-208.
- [9] SHAO M W, LEUNG Y, WU W Z. Rule acquisition and complexity reduction in formal decision contexts[J]. International Journal of Approximate Reasoning, 2014, 55(1): 259-274.
- [10] LI J H, MEI C L, LV Y J. A heuristic knowledge-reduction method for decision formal contexts[J]. Computers and Mathematics with Applications, 2011, 61(4): 1096-1106.
- [11] LI J H, MEI C L, CHERUKURI A K, et al. On rule acquisition in decision formal contexts[J]. International Journal of Machine Learning and Cybernetics, 2013, 4(6): 721-731.
- [12] WU W Z, LEUNG Y, MI J S. Granular computing and knowledge reduction in formal contexts[J]. IEEE Transactions on Knowledge and Data Engineering, 2009, 21(10): 1461-1474.
- [13] 张文修, 仇国芳. 基于粗糙集的不确定性决策[M]. 北京: 清华大学出版社, 2005.
- [14] QU K S, ZHAI Y H, LIANG J Y, et al. Study of decision implications based on formal concept analysis[J]. International Journal of General Systems, 2007, 36(2): 147-156.
- [15] ZHI H L. Extended model of formal concept analysis oriented for heterogeneous data analysis[J]. Acta Electronica Sinica, 2013, 41(12): 2451-2455. (in Chinese)
智慧来. 面向异构数据分析的形式概念分析的扩展模型[J]. 电子学报, 2013, 41(12): 2451-2455.

(上接第 48 页)

- [9] LINGRAS P, CHEN M, MIAO D. Rough cluster quality index based on decision theory[J]. IEEE Transactions on Knowledge and Data Engineering, 2009, 21(7): 1014-1026.
- [10] YU H, LIU Z, WANG G. Automatically determining the number of clusters using decision-theoretic rough set[C]// International Conference on Rough Sets and Knowledge Technology. Banff, Canada, 2011: 504-513.
- [11] MARTIN D, FOWLKES C, TAL D, et al. A database of human segmented natural images and its application to evaluating segmentation algorithms and measuring ecological statistics[C]// Eighth IEEE International Conference on Computer Vision, 2001 (ICCV 2001). IEEE, 2001: 416-423.
- [12] YU Z, AU O C, ZOU R, et al. An adaptive unsupervised approach toward pixel clustering and color image segmentation[J]. Pattern Recognition, 2010, 43(5): 1889-1906.
- [13] BEZDEK J C. Cluster validity with fuzzy sets[J]. Journal of Cybernetics, 1974, 3(3): 58-73.
- [14] BEZDEK J C. Mathematical models for systematics and taxonomy[C]// Proceedings of Eighth International Conference on Numerical Taxonomy. 1975: 143-166.