

学校编码: 10384  
学 号: 200231046

分类号 \_\_\_\_\_ 密级 \_\_\_\_\_  
UDC \_\_\_\_\_

厦门大学  
硕士 学位 论文

基 于 模 糊 聚 类 与 分 类 规 则 的  
数 据 挖 掘 技 术 研 究

**Research on Data Mining Techniques Based on  
Fuzzy Clustering and Classification**

杨 剑 雄

指导教师姓名: 李 茂 青 教 授  
专业 名 称: 系 统 工 程  
论文提交日期: 2005 年 6 月  
论文答辩时间: 2005 年 月  
学位授予日期: 2005 年 月

答 辩 委 员 会 主 席: \_\_\_\_\_  
评 阅 人: \_\_\_\_\_

2005 年 6 月

---

## 厦门大学学位论文原创性声明

兹呈交的学位论文是本人在导师指导下独立完成的研究成果。本人在  
论文写作中参考的其他个人或集体的研究成果均在文中以明确方式标明。  
本人依法享有和承担由此论文而产生的权力和责任。

声明人（签名）：

年 月 日

## 摘要

数据挖掘是利用一种或多种计算机学习技术，从数据库的数据中自动分析并提取有用知识的处理过程。作为一门新兴的交叉性学科，它以数据库技术作为基础，将统计学、逻辑学、机器学习、神经网络、模糊数学、可视化计算等相关领域的成果综合在一起。数据挖掘技术已经成为目前的研究热点。而聚类分析和模式分类作为数据挖掘技术的两个重要课题，在模式识别、图像分割、图形识别等许多方面应用十分广泛。

本文认真研究了数据挖掘技术的基本理论和一般方法，聚类分析和模式分类的一般过程、主要方法以及目前国内的研究现状，并对此进行认真的归纳和总结。在此基础上，对模糊  $c$ -均值算法和基于超盒表示的分类方法进行深入的研究和讨论。

模糊  $c$ -均值算法是目前理论基础较完善、应用较广泛和研究较充分的一种模糊聚类算法，但是它仍存在一些薄弱环节。特别是基于 Mahalanobis 距离的 FCM 算法，收敛速度较慢。本文提出一种改进的算法，通过对算法迭代中聚类中心的移动轨迹图的分析和使用 iris 数据集进行的 200 次实验证明，改进的算法可有效地加快算法的收敛速度，并提高识别率。

由于超盒在规则表示上的方便性，近年来，不少学者开始研究基于超盒表示的分类方法。本文结合模糊聚类提出一种基于模糊聚类的快速分类器算法，构造方法简单、易于理解；针对使用递增学习方法进行训练数据的超盒分类器的共同缺点——分类结果受训练样本排列次序影响，本文提出另一种结合 FCM 的二阶段分类器构建方法，通过与 FMMC 算法的实验对比证明该算法确实有效。将上述两种算法应用于模式识别中取得了一定的效果。

**关键词：**数据挖掘；模糊  $c$ -均值；超盒分类器。

厦门大学博硕士论文摘要库

## ABSTRACT

Data mining is an information extraction activity whose goal is to discover hidden and useful facts contained in databases. Using a combination of machine learning, statistical analysis, modeling techniques and database technology, data mining finds patterns and subtle relationships in data and infers rules that allow the prediction of future results. Typical applications include market segmentation, customer profiling, fraud detection, evaluation of retail promotions, and credit risk analysis.

As a new subject, data mining has drawn the attention of many researchers. Clustering and classification, two of the most important branches of data mining, have been used in many fields such as pattern recognition, classification and recognition of images, etc...

In this paper, the author summarized the theory, methods and techniques used in data mining, as well as current status of research in this field around the world.

The main focus of the paper was on Fuzzy  $c$ -means (FCM) algorithm and hyperbox-based classifier methods.

While theoretically well developed and widely used, there are still some drawbacks in FCM, for example, the huge execution time for Mahalanobis-distance-based FCM. In this paper, the author proposed a modified algorithm(MFCM) to improve the convergent time of FCM-based hyperelliptical clusters and tested it on Fisher's famous iris data set. The outcomes showed that MFCM can decrease the iterative times and quicken the convergent speed effectively.

In recent years, more and more people focus their research on hyperbox-based classifiers due to the easy extraction of IF-THEN rules from hyperbox. The author also proposed two algorithms to generate hyperbox-based classifiers and tested them on a heart disease data set. The outcomes proved that the algorithms were effective and efficient.

**Keywords:** Data mining; Fuzzy  $c$ -means; Hyperbox-based classifier;

厦门大学博硕士论文摘要库

## 目 录

<b>第一章 绪论 .....</b>	<b>1</b>
<b>1.1 数据挖掘概述 .....</b>	<b>1</b>
1.1.1 引言 .....	1
1.1.2 数据挖掘的定义 .....	2
1.1.3 数据挖掘的功能 .....	2
1.1.4 数据挖掘方法 .....	4
<b>1.2 基于模糊聚类与分类规则的数据挖掘技术研究现状 .....</b>	<b>6</b>
1.2.1 模糊聚类 .....	6
1.2.2 模式分类 .....	8
<b>1.3 本文的研究内容与结构安排 .....</b>	<b>9</b>
<b>第二章 基于分类规则的数据挖掘 .....</b>	<b>11</b>
<b>2.1 分类和分类的一般过程 .....</b>	<b>11</b>
<b>2.2 常用的分类方法 .....</b>	<b>11</b>
2.2.1 线性判别分类法 .....	11
2.2.2 贝叶斯分类法 .....	12
2.2.3 决策树分类法 .....	15
<b>2.3 基于超盒表示的分类器 .....</b>	<b>18</b>
2.3.1 引言 .....	18
2.3.2 FMMC 算法 .....	18
2.3.3 NGE 算法 .....	22
2.3.4 Fuzzy ARTMAP 算法 .....	23
2.3.5 ALFC 算法 .....	24

---

2.3.6 小结 .....	25
<b>第三章 面向聚类的数据挖掘 .....</b>	<b>26</b>
<b>3.1 引言 .....</b>	<b>26</b>
<b>3.2 聚类分析方法的分类.....</b>	<b>26</b>
3.2.1 划分方法（Partitioning Methods） .....	26
3.2.2 层次法（Hierarchical Methods） .....	27
3.2.3 基于密度的方法（Density-Based Methods） .....	27
3.2.4 基于网格的方法（Grid-Based Methods） .....	28
3.2.5 基于模型的方法（Model-Based Methods） .....	28
<b>3.3 模糊聚类分析.....</b>	<b>28</b>
3.3.1 模糊集定义与隶属函数 .....	28
3.3.2 模糊 $c$ -均值算法.....	29
<b>第四章 改进的模糊 <math>C</math>-均值算法 .....</b>	<b>37</b>
<b>4.1 引言 .....</b>	<b>37</b>
<b>4.2 算法描述 .....</b>	<b>38</b>
<b>4.3 改进的 FCM 算法（MFCM）的步骤.....</b>	<b>40</b>
<b>4.4 实验过程与讨论 .....</b>	<b>41</b>
<b>4.5 MFCM 算法的性能分析 .....</b>	<b>43</b>
<b>4.6 小结 .....</b>	<b>46</b>
<b>第五章 基于模糊聚类和分类规则的模式识别方法研究 .....</b>	<b>47</b>
<b>5.1 基于特征空间的模糊聚类算法及其在模式识别中的应用 .....</b>	<b>47</b>
5.1.1 引言 .....	47
5.1.2 基于特征空间的模糊聚类（FSFC）算法描述 .....	47

5.1.3 FSFC 算法过程 .....	48
5.1.4 实验过程与讨论 .....	49
<b>5.2 结合模糊聚类的二阶段超盒分类器的构建及应用 .....</b>	<b>51</b>
5.2.1 结合模糊聚类的二阶段超盒分类器（HTSC）算法描述 .....	51
5.2.2 实验过程与结果 .....	55
5.2.3 HTSC 算法的分析与改进 .....	56
<b>5.3 小结 .....</b>	<b>58</b>
<b>第六章 结束语 .....</b>	<b>59</b>
6.1 本文的主要工作 .....	59
6.2 今后进一步的研究工作 .....	59
<b>参考文献 .....</b>	<b>61</b>
<b>致 谢 .....</b>	<b>64</b>
<b>附录 1 .....</b>	<b>65</b>
<b>附录 2 .....</b>	<b>66</b>

厦门大学博硕士论文摘要库

## Contents

<b>Chapter 1 Introduction.....</b>	<b>1</b>
<b>1.1 Data Mining .....</b>	<b>1</b>
1.1.1 Introduction .....	1
1.1.2 Definition of Data Mining.....	2
1.1.3 Data Mining Functionalities.....	2
1.1.4 Methods of Data Mining .....	4
<b>1.2 Current Status of Research on Fuzzy Clustering &amp; Classification-Based Data Mining .....</b>	<b>6</b>
1.2.1 Development of Fuzzy Clustering .....	6
1.2.2 Development of Classification.....	8
<b>1.3 Table of Contents of this paper .....</b>	<b>9</b>
<b>CHAPTER 2 CLASSIFICATION.....</b>	<b>11</b>
<b>2.1 What is Classification?.....</b>	<b>11</b>
<b>2.2 Classification Methods.....</b>	<b>11</b>
2.2.1 Classification by Linear Discriminant.....	11
2.2.2 Bayesian Classification .....	12
2.2.3 Classification by Decision Tree Induction .....	15
<b>2.3 Hyperbox-based Classifier .....</b>	<b>18</b>
2.3.1 Introduction .....	18
2.3.2 FMMC Algorithm .....	18
2.3.3 NGE Algorithm .....	22
2.3.4 Fuzzy ARTMAP Algorithm.....	23
2.3.5 ALFC Algorithm .....	24
2.3.6 Summary .....	25

---

**CHAPTER 3 CLUSTERING ANALYSIS ..... 26**

<b>3.1 Introduction .....</b>	<b>26</b>
<b>3.2 A Categorization of Major Clustering Methods .....</b>	<b>26</b>
3.2.1 Partitioning Methods .....	26
3.2.2 Hierarchical Methods .....	27
3.2.3 Density-Based Methods .....	27
3.2.4 Grid-Based Methods .....	28
3.2.5 Model-Based Methods .....	28
<b>3.3 Fuzzy Clustering Analysis .....</b>	<b>28</b>
3.3.1 Fuzzy Set and Membership Function.....	28
3.3.2 Fuzzy C-Means Algorithm.....	29

**CHAPTER 4 A MODIFIED FUZZY C-MEANS(MFCM)****ALGORITHM..... 37**

<b>4.1 Introduction .....</b>	<b>37</b>
<b>4.2 MFCM Description .....</b>	<b>38</b>
<b>4.3 MFCM Algorithm .....</b>	<b>40</b>
<b>4.4 Experimental Results .....</b>	<b>41</b>
<b>4.5 Performance Analysis of MFCM .....</b>	<b>43</b>
<b>4.6 Conclusion.....</b>	<b>46</b>

**CHAPTER 5 RESEARCH ON PATTERN RECOGNITION  
BASED ON FUZZY CLUSTERING AND  
CLASSIFICATION..... 47**

<b>5.1 Feature Space-Based Fuzzy Clustering Algorithm (FSFC) and its application for Pattern Recognition .....</b>	<b>47</b>
5.1.1 Introduction .....	47
5.1.2 FSFC Description.....	47

5.1.3 FSFC Algorithm .....	48
5.1.4 Experimental Results.....	49
<b>5.2 A Hyperbox-Based Two-Stage Classifier(HTSC) Algorithm.....</b>	<b>51</b>
5.2.1 HTSC Description.....	51
5.2.2 Experimental Results.....	55
5.2.3 Analysis and Improvement of HTSC .....	56
<b>5.3 Conclusion.....</b>	<b>58</b>
<b>CHAPTER 6 SUM-UP AND FURTHER RESEARCH .....</b>	<b>59</b>
<b>REFERENCES .....</b>	<b>61</b>
<b>ACKNOWLEDGEMENT.....</b>	<b>64</b>
<b>APPENDIX 1 .....</b>	<b>65</b>
<b>APPENDIX 2 .....</b>	<b>66</b>

厦门大学博硕士论文摘要库

# 第一章 绪论

## 1.1 数据挖掘概述

### 1.1.1 引言

在当今的信息社会时代中，信息已经成为社会的三大重要资源之一，如何有效地管理和利用这些信息资源成为人们关注的焦点与技术发展的趋势。大量信息在给人们带来方便的同时，也带来了一系列问题：如何从大量信息中提取和选择所需数据、信息的真伪辨析、信息组织形式的一致化等等。计算机的快速发展给信息处理带来了巨大的方便，但随着人们对信息的需求和要求的不断增长，对数据处理技术的要求也不断的提高，数据挖掘（Data Mining, DM）即是由此而发展起来的一项新技术。

数据挖掘出现于 20 世纪 80 年代末，最早是以知识发现——从数据库中发现知识（Knowledge Discovery in Database, KDD）的研究起步的。数据挖掘是 KDD 中的一个重要步骤，但在通常的应用中，并不区分二者的概念。KDD 一词首先出现在 1989 年的人工智能国际会议上，之后 KDD 的研究得到迅速发展，自 1995 年起，知识发现与数据挖掘国际学术会议每年均召开一届，此外还有相关的地区性国际大会定期或不定期召开。目前，数据库中的知识发现和数据挖掘技术已成为研究热点和焦点，在国外已达到一定的水平并投入应用领域中，商品化的 KDD 软件工具已投放到市场，如 IBM 公司的 IBM DB2 Intelligent Miner、美国 Business Objects 公司的 Business Miner 等。国内从事数据挖掘的研究起步较晚，但近年来许多高校、科研院所在这一领域的研究也取得一定的成绩。

数据挖掘的应用领域非常广泛，如生物医学、零售、金融、气象、电

---

子商务等等.

### 1.1.2 数据挖掘的定义

对数据挖掘的定义，不同的学者有不同的表述形式，归结起来，数据挖掘是一个处理过程，是从大量的、不完全的、有噪声的、模糊的、随机的实际应用数据中，提取隐含在其中的，人们事先不知道的，但又是潜在的有用的并且最终可理解的信息和知识的过程。<sup>[1]</sup>

### 1.1.3 数据挖掘的功能

数据挖掘的任务是从数据集中发现模式，模式有许多种，按功能可分为两大类：预测型（Predictive）和描述型（Descriptive），而在实际应用中，往往根据模式的实际作用进行细分，如分类、聚类、回归等。数据挖掘所处理的数据类型也很丰富，包括数值数据、文本数据、关系数据库、Web 网页等等。一般来说，数据挖掘的功能与被挖掘的数据类型有关，某些功能只能应用于特定数据类型中，而某些功能则可以应用于多种数据类型中。数据挖掘的功能具体包括以下内容<sup>[2]</sup>.

- 概念描述

概念描述（Concept Description，或称为类描述 Class Description）是通过对与某类对象关联数据的汇总、分析和比较，对此类对象的内涵进行描述，并概括这类对象的有关特征。这种描述是汇总的、简洁的和精确的，并且能够对数据分析起到重要作用。概念描述可以通过数据特征化（data characterization）和数据区分（data discrimination）两种方法获得。前者是对目标类数据的一般特征或特性的汇总，生成的描述一般是该类中所有对象的共性；后者是将目标类数据的一般特性与一个或多个对比类数据的一般特性进行比较，生成的描述一般是目标类和对比类中对象的共性。

Degree papers are in the "[Xiamen University Electronic Theses and Dissertations Database](#)". Full texts are available in the following ways:

1. If your library is a CALIS member libraries, please log on <http://etd.calis.edu.cn/> and submit requests online, or consult the interlibrary loan department in your library.
2. For users of non-CALIS member libraries, please mail to [etd@xmu.edu.cn](mailto:etd@xmu.edu.cn) for delivery details.

厦门大学博硕士论文摘要库