



MOHFENLEI JIQIZAI

GUANGPUXINXICHULIZHONGDE YINGYONG

# 模糊分类及其在

# 光谱信息处理中的应用

武小红 武斌 ◎著

东南大学出版社  
SOUTHEAST UNIVERSITY PRESS

# 模糊分类及其在光谱 信息处理中的应用

武小红 武斌 著



SE 东南大学出版社  
SOUTHEAST UNIVERSITY PRESS

·南京·

## 内 容 简 介

本书主要研究模糊分类中的模糊聚类,模糊特征提取和模糊分类器以及它们在农产品/食品的近红外光谱信息处理中的应用。模糊聚类算法涉及模糊 C-均值聚类、联合模糊 C-均值聚类、利用核方法和新的非欧氏距离改进一些模糊聚类、一种改进的可能模糊 C-均值聚类算法等。模糊特征提取涉及模糊线性判别分析、核模糊主成分分析、核模糊判别分析、模糊非相关判别分析等。模糊分类器涉及模糊 K 近邻,核模糊 K 近邻。用模糊聚类算法对苹果近红外光谱、生菜近红外光谱、茶叶中红外光谱进行分类。用模糊线性判别分析和模糊非相关判别分析提取近红外光谱的鉴别信息。

本书可为研究模糊分类、模式识别和近红外光谱信息处理方向的科研工作者提供借鉴和参考,也可为从事农产品、食品近红外光谱信息的分析和处理的研究生和教师提供参考。

## 图书在版编目(CIP)数据

模糊分类及其在光谱信息处理中的应用 / 武小红,  
武斌著. —南京:东南大学出版社,2017. 10

ISBN 978 - 7 - 5641 - 7182 - 7

I . ①模… II . ①武… ②武… III . ①模糊分  
类—应用—红外光谱—信息处理 IV . ①TP274

中国版本图书馆 CIP 数据核字(2017)第 243587 号

## 模糊分类及其在光谱信息处理中的应用

---

出版发行 东南大学出版社

出版人 江建中

社 址 南京市四牌楼 2 号

邮 编 210096

---

经 销 全国各地新华书店

印 刷 虎彩印艺股份有限公司

开 本 700 mm×1000 mm 1/16

印 张 10.75

字 数 211 千字

版 次 2017 年 10 月第 1 版

印 次 2017 年 10 月第 1 次印刷

书 号 ISBN 978 - 7 - 5641 - 7182 - 7

定 价 45.00 元

---

(本社图书若有印装质量问题,请直接与营销部联系。电话:025-83791830)

# 前　言

模糊分类是模式识别中一个重要的分支,它是模糊数学在信息科学中的应用之一。当精确数学方法无法处理自然和社会中的模糊事物时,模糊数学随着科技发展的需要诞生了。当我们需要对一些模糊的事物进行分类时,模糊分类比传统分类方法更能够体现事物的不确定性。在模糊分类算法中模糊聚类应用得最广泛,它广泛应用于汉字字符识别、语音识别、图像处理和雷达目标识别等等。但是,涉及应用模糊分类处理近/中红外光谱信息方面的相关文献比较少,而这正是本书的主要论述内容。

本书主要论述了模糊分类中的一些算法,尤其是模糊聚类分析。同时,将模糊分类应用到农产品/食品的近红外/中红外光谱信息处理中。如何提高聚类准确率,降低噪声的影响,减少聚类时间是模糊聚类分析时要解决的主要问题,也是本书主要的论述内容。

本书主要研究模糊分类中的三大部分:模糊聚类,模糊特征提取和模糊分类器。在分析现有的几种模糊聚类基础上提出了一些新的模糊聚类算法以提高聚类的准确性,减少聚类时间和提高聚类性能,同时结合核方法和新的距离测度改进模糊聚类算法,研究苹果近红外光谱、生菜近红外光谱和茶叶中红外光谱的模糊聚类分析。本书的主要内容包括:第1章介绍了模糊分类的发展概况和基本理论知识;第2章论述了几种模糊混合聚类模型;第3章应用核方法改进现有的几种模糊聚类算法;第4章论述了基于非欧氏距离的模糊聚类算法;第5章论述了模糊鉴别信息提取算法及基于核的模糊鉴别信息提取算法和基于核的模糊K-近邻法;第6章论述了两种新的GK聚类算法以及它们在茶叶红外光谱分类中的应用;第7章论述了苹果近红外光谱的模糊聚类分析;第8章论述了模糊鉴别C均值聚类,模糊鉴别学习矢量量化和一种广义噪声聚类以及它们在茶叶红外光谱分类中的应用;第9章论述了两种新的模糊K调和均值聚类以及它们在光谱分类中的应用;第10章论述了四种模糊学习矢量量化模型以及它们在光谱分类中的应用。

本书主要由江苏大学武小红副教授(撰写了15.1万字)和滁州职业技术学院武斌副教授(撰写了6万字)完成。在本书写作过程中,孙俊教授、李敏教授给予了大力支持并提供了宝贵的建议和帮助。本书的出版获得了安徽省高等教育振兴计划人才项目“高校优秀青年人才支持计划”(皖教秘人〔2014〕181号)和江苏高校优势学科建设工程资助项目PAPD的资助。在此表示衷心的感谢。

由于作者业务水平和时间所限,书中难免存在错误和不当之处,敬请读者批评指正。

# 目 录

<b>1</b>	<b>绪论</b>	( 1 )
1.1	模糊分类概述	( 1 )
1.2	模糊分类的发展概况	( 2 )
1.3	模糊分类的基本理论简介	( 4 )
1.3.1	模糊集合	( 4 )
1.3.2	基于目标函数的模糊聚类	( 5 )
1.3.3	模糊判别分析	( 7 )
1.4	本章小结	( 8 )
	参考文献	( 8 )
<b>2</b>	<b>模糊混合聚类模型</b>	( 11 )
2.1	模糊聚类概述	( 11 )
2.2	联合模糊 C-均值聚类模型	( 12 )
2.2.1	MPCM 模型	( 12 )
2.2.2	AFCM 模型	( 13 )
2.2.3	MPCM 和 AFCM 的对比实验	( 17 )
2.3	一种改进的可能模糊聚类算法	( 20 )
2.3.1	PCA 算法及其存在的问题	( 20 )
2.3.2	改进的 PCM 与改进的 PFCM	( 23 )
2.3.3	实验结果	( 27 )
2.4	基于聚类中心分离的可能模糊聚类模型	( 31 )
2.4.1	基于聚类中心分离的模糊 C-均值聚类(FCM_CCS)	( 31 )
2.4.2	基于聚类中心分离的可能聚类(PCM_CCS)	( 32 )
2.4.3	基于聚类中心分离的可能模糊聚类(PFCM_CCS)	( 33 )
2.4.4	实验结果	( 34 )
2.5	一种混合可能聚类算法	( 38 )
2.5.1	算法描述	( 38 )
2.5.2	实验结果	( 39 )

2.6 联合模糊熵聚类 .....	( 40 )
2.6.1 算法描述 .....	( 40 )
2.6.2 实验结果 .....	( 42 )
2.7 本章小结 .....	( 43 )
参考文献 .....	( 44 )

### **3 基于核的模糊聚类 .....** ( 46 )

3.1 核模糊聚类概述 .....	( 46 )
3.2 基于核的修正可能 C-均值聚类 .....	( 46 )
3.2.1 算法描述 .....	( 46 )
3.2.2 实验结果 .....	( 48 )
3.3 基于核的广义噪声聚类算法 .....	( 50 )
3.3.1 GNC 算法 .....	( 50 )
3.3.2 KGNC 算法 .....	( 51 )
3.3.3 实验结果 .....	( 53 )
3.4 基于核的可能模糊 C-均值聚类 .....	( 55 )
3.4.1 算法描述 .....	( 55 )
3.4.2 实验结果 .....	( 56 )
3.5 基于核的聚类中心分离的模糊 C-均值聚类 .....	( 57 )
3.5.1 算法描述 .....	( 57 )
3.5.2 实验结果 .....	( 58 )
3.6 基于核的类间分离聚类 .....	( 58 )
3.6.1 算法描述 .....	( 58 )
3.6.2 实验结果 .....	( 60 )
3.7 本章小结 .....	( 60 )
参考文献 .....	( 60 )

### **4 基于非欧氏距离的模糊聚类算法 .....** ( 63 )

4.1 引言 .....	( 63 )
4.2 一种新的非欧氏距离 .....	( 64 )
4.3 基于非欧氏距离可能模糊 C-均值聚类算法 .....	( 65 )
4.3.1 可能模糊 C-均值聚类算法 .....	( 65 )
4.3.2 APFCM 算法 .....	( 66 )
4.3.3 实验结果 .....	( 67 )

4.4 基于非欧氏距离可能聚类算法 .....	(69)
4.4.1 IPCM 算法 .....	(69)
4.4.2 AIPCM 算法 .....	(70)
4.4.3 实验结果 .....	(71)
4.5 本章小结 .....	(73)
参考文献 .....	(74)
<b>5 基于核的模糊鉴别信息提取及分类 .....</b>	<b>(75)</b>
5.1 引言 .....	(75)
5.2 基于核的模糊判别分析(KFDA) .....	(76)
5.2.1 KFDA 算法 .....	(76)
5.2.2 实验结果 .....	(78)
5.3 模糊主元分析及其核模型 .....	(79)
5.3.1 模糊主元分析 .....	(79)
5.3.2 基于核的模糊主元分析 .....	(79)
5.3.3 实验结果 .....	(81)
5.4 模糊非相关判别转换(FUDT)及其核模型 .....	(82)
5.4.1 非相关判别转换(UDT) .....	(82)
5.4.2 模糊非相关判别转换(FUDT) .....	(84)
5.4.3 基于核的模糊非相关判别分析 .....	(86)
5.4.4 苹果近红外光谱的线性和非线性鉴别信息提取实验 .....	(88)
5.5 基于核的模糊 K-近邻法 .....	(90)
5.5.1 模糊 K-近邻法 .....	(90)
5.5.2 基于核的模糊 K-近邻法 .....	(91)
5.5.3 实验结果 .....	(92)
5.6 本章小结 .....	(94)
参考文献 .....	(94)
<b>6 基于模糊协方差矩阵聚类的茶叶红外光谱分类 .....</b>	<b>(97)</b>
6.1 一种混合 GK 聚类 .....	(98)
6.1.1 算法描述 .....	(98)
6.1.2 实验结果 .....	(99)
6.2 模糊协方差矩阵的可能模糊聚类 .....	(101)
6.2.1 算法描述 .....	(102)

6.2.2 实验结果 .....	(103)
6.3 本章小结 .....	(104)
参考文献 .....	(105)

**7 苹果近红外光谱的模糊聚类分析 .....** (107)

7.1 苹果近红外光谱检测研究 .....	(107)
7.1.1 国外研究进展 .....	(108)
7.1.2 国内研究进展 .....	(108)
7.2 苹果近红外光谱模糊聚类 .....	(110)
7.2.1 GK 和 GG 模糊聚类 .....	(110)
7.2.2 一种混合模糊类间分离聚类 .....	(111)
7.2.3 实验结果 .....	(114)
7.3 一种快速联合模糊 C-均值聚类 .....	(117)
7.3.1 FAFCM 聚类 .....	(117)
7.3.2 实验结果 .....	(118)
7.4 本章小结 .....	(120)
参考文献 .....	(121)

**8 茶叶傅里叶红外光谱模糊聚类分析 .....** (125)

8.1 模糊鉴别 C-均值聚类 .....	(125)
8.1.1 算法描述 .....	(125)
8.1.2 实验结果 .....	(126)
8.2 模糊鉴别学习矢量量化 .....	(129)
8.2.1 算法描述 .....	(129)
8.2.2 实验结果 .....	(130)
8.3 一种广义噪声聚类 .....	(131)
8.3.1 算法描述 .....	(131)
8.3.2 实验结果 .....	(132)
8.4 本章小结 .....	(134)
参考文献 .....	(135)

**9 模糊 K 调和均值聚类的近/中红外光谱分类 .....** (137)

9.1 K 调和均值聚类 .....	(137)
9.2 广义模糊 K 调和均值聚类的近红外光谱生菜储藏时间鉴别 .....	(138)
9.2.1 算法描述 .....	(138)

---

9.2.2 实验结果	(140)
9.3 一种混合模糊 K 调和均值聚类	(142)
9.3.1 算法描述	(144)
9.3.2 实验结果	(144)
9.4 本章小结	(145)
参考文献	(146)
<b>10 模糊学习矢量量化模型</b>	<b>(148)</b>
10.1 可能模糊学习矢量量化	(148)
10.1.1 算法描述	(148)
10.1.2 实验结果	(149)
10.2 无监督可能模糊学习矢量量化的近红外光谱生菜品种鉴别	(150)
10.2.1 算法描述	(151)
10.2.2 实验结果	(152)
10.3 一种基于优化的模糊学习矢量量化的苹果分类	(155)
10.3.1 算法描述	(155)
10.3.2 实验结果	(156)
10.4 联合模糊学习矢量量化	(156)
10.4.1 算法描述	(156)
10.4.2 实验结果	(157)
10.5 本章小结	(157)
参考文献	(158)

## 1.1 模糊分类概述

科学的目的是发现事物本身的内在规律,也就是去识别事物的“模式”。模式识别是随着科学技术的进步在20世纪60年代迅速发展起来的一门新兴学科。模式识别理论和方法的不断完善和进步推动了相关学科(例如人工智能)的发展,它已经广泛应用于信息处理、人工智能、医疗诊断、信息安全、国防等众多领域。模糊分类是模式识别的一个重要分支,它的诞生是模糊数学应用在信息科学领域的实例。1965年,美籍伊朗著名控制论专家L.A.Zadeh发表了开创性的论文《模糊集合》,正式提出了多值集合理论,创立了模糊数学,从而提供了一套严格的数学方法,用来描述这种带有模糊不确定性的现象和事物。模糊数学决不是模模糊糊的数学,现实世界中往往不需要绝对的精确,一定的不准确程度是可以接受的,当精确的数学方法无法处理现实世界中一些不确定(模糊)事物时,模糊数学往往能够很好地处理。模糊信息处理就是利用模糊数学这一工具,来处理带有模糊不确定性的信息。由于现实世界中许多现象中含有模糊性信息,因此模糊分类技术自从其诞生开始就表现出了广阔的应用前景。

近年来模糊分类理论得到了国内外学者的高度重视,取得了很好的发展。在一些具体的识别应用中,模糊分类也具有较好的效果,比如图像处理、汉字字符识别、语音识别和雷达目标识别等等。随着科学技术的发展,模糊数学的不断完善和进步,必然对模糊分类技术带来更多的挑战,因此必须努力丰富和完善模糊分类理论和实际应用技术以满足人们生产和生活的需要。

模糊分类是一种机器学习(Machine Learning)算法,按学习过程中有无教师参与,分为有监督学习和无监督学习。

进行有监督学习分类时,已知样本的类别和类别属性,先让学习算法或分类器对有类别标记的样本进行学习或训练,然后用训练好的分类器对待分类样本进行分类。这种分类需要有足够的学习样本以便得到充分的先验知识。

进行无监督学习分类时,在没有任何的先验知识,没有教师参与,也就是分类器没有经过学习知识的情况下进行未知样本的分类。

本书主要探讨模糊分类中的模糊聚类、模糊特征提取和模糊分类器。聚类(Clustering)就是寻找数据之间的相互联系,按照相似性原则进行分类,就是“物以类聚”。传统的硬聚类算法使某个数据点具有严格的隶属关系,隶属度要么是1要么是0。因而硬聚类并不能很好地反映数据点对类的实际隶属关系,在现实世界中这种隶属关系往往是模糊的。Zadeh的模糊集理论为处理这种模糊关系提供了有力的分析工具,用模糊的方法来处理聚类问题则称之为模糊聚类。目前,大多数模糊聚类算法是无监督的学习算法,也有少数是带有监督的模糊聚类算法,主要研究无监督的模糊聚类算法。判别分析(Discriminant Analysis)是一种根据观测变量判断样本如何分类的多元统计方法,对预测变量的线性组合产生一系列判别函数,利用判别函数决定待分类样本的类别归属。常见的判别分析有贝叶斯(Bayes)判别、Fisher 线性判别、Foley-Sammon 判别、非相关线性(Uncorrelated Discriminant)判别等。将模糊集理论应用到判别分析就得到模糊判别分析。

## 1.2 模糊分类的发展概况

### 1) 模糊聚类的形成概况

1966 年 Bellman, Kalaba 和 Zadeh 最早将模糊集思想引入到聚类算法中,不久 Wee, Flake 和 Turner, Gitman 和 Levine 做了一定的尝试。1970 年 Ruspini 首次系统地提出模糊聚类算法,该算法建立在最小化一个模糊目标函数的基础上,这对以后的模糊聚类算法有很大的启发作用。1973 年 Duda 和 Hart 按照类内平方误差(WGSS: within-group sum of squared errors)最小准则提出 ISODATA 算法(即硬 C-均值聚类算法)。Dunn 将 ISODATA 算法推广到加权的 ISODATA 算法,后由 Bezdek 推广到模糊 C-均值聚类(FCM)算法。

### 2) 模糊聚类中距离测度的研究概况

1979 年 Gustafson 和 Kessel 引入模糊协方差矩阵,提出了基于协方差矩阵加权的距离度量(Mahalonobis 距离)代替 FCM 的 Euclidean 距离得到 GK 聚类算法,GK 聚类算法适合处理同一个数据集中包含不同拓扑结构的数据。欧氏距离只能检测球状数据,为拓宽模糊聚类的适用范围,1981 年 Bezdek 提出了模糊 C 族算法,采用样本点到族原型的距离度量样本与各类原型的相似性,当族维数取 0,1 或 2 以上的值时,分别适用于椭球状、线状、平面状、超平面状的数据结构。1989 年 Gath 和 Geva 将最大相关原则引入距离的计算中,提出另一种距离度量方式。1991 年 Bobrowski 和 Bezdek 提出了  $L_1$  和  $L_\infty$  范数(即 Hamming 和 Maximum 距离)下的模糊聚类算法,发现在许多情况下它们比常用的欧氏范数  $L_2$  能获得更好的结果,建议在聚类分析中要选择合适的距离函数。1992 年 Dave 提出了模糊 C

壳算法,采用点到圆壳的垂直距离度量方式,模糊 C 壳算法的一些修改算法也相继提出。2001 年,为了更好地处理彩色图像 Ozdemir 提出了 ICS 聚类算法,在模糊 C-均值算法的目标函数中加入一个分裂项。2002 年 Timm 和 Kruse 为了避免产生一致的类中心,在 PCM 算法的目标函数中也加入类中心的排斥提出一种新 PCM 算法。2002 年 Wu 和 Yang 在鲁棒统计观点和影响函数基础上提出了一种新的非欧氏距离以替代 FCM 和 PCM 中的欧氏距离。此外,针对关系型数据和模糊数,分别提出了关系型模糊 C-均值聚类算法和模糊数型模糊聚类算法。借助不同的距离度量方式即可构造不同的目标函数,形成不同的分类标准,但是如何选取合适的度量准则目前尚缺乏理论指导。

### 3) 模糊聚类对含噪声数据处理的研究概况

1993 年为了解决 FCM 对噪声敏感的问题,Krishnapuram 和 Keller 放弃了 FCM 的可能性约束条件,构造了一个新的目标函数,提出了可能 C-均值聚类(PCM)。PCM 能够聚类包含噪声或野值点的数据,PCM 使噪声数据具有很小的隶属度值,因而噪声对聚类的影响可以忽略。但是 Barni、Cappellini 和 Mecocci 通过实验发现 PCM 对初始化条件很敏感,常常会导致一致性聚类结果。Pal 等为解决 PCM 对初始聚类中心很敏感,常常会导致一致性聚类结果这个问题,提出了混合 C-均值(Mixed C-Means)聚类算法。但是如果数据量大时混合 C-均值聚类得出的隶属度将很小,这是由于混合 C-均值聚类引入了新的约束条件所导致的。2004 年 Pal、Bezdek 等提出了可能性模糊 C-均值聚类(PFCM)很好地解决了 PCM 的问题。Dave 从另外一个角度出发解决了 FCM 对噪声敏感的问题。将噪声点看成单独的一个类,提出了噪声聚类算法和广义噪声聚类算法以处理噪声数据。

### 4) 模糊聚类有效性的研究

模糊聚类有效性是探讨模糊聚类算法对数据集聚类的结果是否合理,对数据集的划分是否反映了未知数据结构的真实情况。这样就需要使用一些有效性准则。第一个模糊有效性准则函数由 Bezdek 为 FCM 设计的划分系数(Partition Coefficient)。Bezdek 的另一个模糊有效性准则函数称为划分熵(Partition Entropy)。划分系数和划分熵只和模糊隶属度有关而没有考虑到数据集的数据结构因而具有局限性。Xie 和 Beni 提出的 Xie-Beni 模糊有效性准则函数以及 Fukayama 和 Sugno 提出的 Fukayama-Sugno 模糊有效性准则函数可以度量数据结构的紧凑性。Pal 和 Bezdek 指出权重指数  $m$  取值不同时对 Fukayama-Sugno 有效性函数影响较大而对 Xie-Beni 有效性函数较小,并建议  $m$  取值范围为 [1.5 2.5]。Gath 和 Geva 在基于超体积和密度的概念下提出三种有效性函数,他们认为好的分类效果应该具有低的模糊超体积。SUN 等提出一种新的有效性函数  $V_{wsj}$  以衡量类间分离度和类内紧凑度,最小化  $V_{wsj}$  得到最好的聚类效果。还有 Zahid、

Bouguessa 等学者提出的一些有效性准则函数。

### 5) 模糊聚类在国内的研究概况

模糊聚类在国内的研究和应用也很广泛。郭桂蓉等较早地利用模糊模式识别进行雷达目标识别研究,庄钊文等进一步推广了模糊模式识别在雷达目标识别领域的研究。伍忠东将核方法应用到模糊 C-均值聚类提出核模糊 C-均值聚类。高新波提出加权模糊 C-均值聚类,对模糊 C-均值聚类算法中加权指数  $m$  进行了研究。于剑提出一种新的 FCM 有效性系数。范九伦等用可能性分布的观点解释划分系数,提出了一个新的有效性函数,其性能优于划分系数。罗建军利用数据点特征权重的概率约束关系和可能分布,提出了分别建立在概率和可能加权特征方式之上的改进可能模糊聚类的两种模型。李洁通过 ReliefF 算法对特征进行加权选择,不仅能够将模糊  $k$ -均值、 $k$ -mode 以及  $k$ -原型算法合而为一,同时使样本的分类效果更好,在加权模糊 C-均值聚类基础上提出基于特征加权的模糊聚类新算法。为了克服 PCM 的缺点,Zhang 和 Leung 在模糊 C-均值聚类(FCM)和可能性 C-均值聚类(PCM)的基础上提出了改进型可能性 C-均值聚类(IPCM)。国内还有很多有关模糊聚类在医学图像分割、图像处理、水文信息处理、汉字字符识别、雷达目标识别等等方面的应用,这里不一一赘述。

### 6) 模糊分类的分类器研究概况

同大量的模糊聚类算法相比,模糊分类中的有监督学习算法明显比较少。1985 年 Keller 等将模糊理论运用到 K-近邻法提出模糊 K-近邻法。1999 年 Chen 等提出模糊线性判别分析并用此算法对化学数据集做实验,实验表明该算法比传统的线性判别分析性能优良。Kwak 和 Pedrycz 提出了模糊 fisherface 分类器对人脸进行分类识别,可以减少因为光照、观察环境和脸部表情等因素的变化而影响分类效果。Lin 和 Wang 提出了模糊支持向量机(Fuzzy Support Vector Machines),将每个数据点用一个隶属度来反映数据的关联程度,能较好地处理噪声数据和野值数据。Tsujinishi 和 Abe 为解决多类不可分区域的分类问题提出了模糊最小平方支持向量机。

### 7) 模糊分类的应用

主要应用在人脸识别、图像处理、文本识别、雷达目标识别和光谱信息处理等领域。

## 1.3 模糊分类的基本理论简介

### 1.3.1 模糊集合

模糊分类是建立在模糊数学中的模糊集合理论基础上的分类算法。

经典集合:对于给定的论域  $X$  和某一属性  $P$ , $X$  中满足属性  $P$  的所有元素组

成的全体叫做集合。

模糊集合：对于给定的论域  $X$ ,  $x$  是其中的元素。对于任意的  $x \in X$ , 给定了如下的映射：

$$X \rightarrow [0,1]; x \mapsto u_A(x) \in [0,1],$$

则称如下的集合为上的模糊集合：

$$\tilde{A} = \{(x | \mu_{\tilde{A}}(x))\}, \forall x \in X$$

称  $\mu_{\tilde{A}}(x)$  为元素  $x$  对  $\tilde{A}$  的隶属度, 它反映了元素  $x$  对  $\tilde{A}$  的隶属程度。

经典集合是只考虑“非此即彼”的集合, 从论域  $X$  中任意取出一个元素  $x$ ,  $x$  要么属于集合  $A$  要么不属于集合  $A$ 。但是在现实世界中很多事物的隶属边界没有那么分明, 边界是不清晰的、模糊的。为了用数学方法来解决这种问题, 1965 年 Zadeh 教授提出了模糊集合的概念。他用隶属度来描述处于中间过渡状态的事物对差异双方所具有的倾向性。隶属度是经典集合中的特征函数的推广。将经典集合的特征函数的值域由  $\{0,1\}$  二值扩展到  $[0,1]$  区间时, 就描述了一个模糊集合。描述模糊集合的工具是隶属度, 如果隶属度  $\mu_{\tilde{A}}(x)$  接近于 1, 则表示  $x$  对于模糊集合  $\tilde{A}$  的隶属程度很高; 反之, 如果隶属度  $\mu_{\tilde{A}}(x)$  接近于 0, 则表示  $x$  对于模糊集合  $\tilde{A}$  的隶属程度很低。

模糊集合的表示方法：

(1) 设论域  $X$  是有限的可数集, 令  $X = \{x_1, x_2, \dots, x_n\}$ ,  $X$  上的任一模糊集  $\tilde{A}$ , 其隶属度为  $\mu_{\tilde{A}}(x_i), i=1, 2, \dots, n$ , 则  $\tilde{A}$  可以表示成

$$\tilde{A} = \mu_{\tilde{A}}(x_1)/x_1 + \mu_{\tilde{A}}(x_2)/x_2 + \dots + \mu_{\tilde{A}}(x_n)/x_n = \sum_{i=1}^n \mu_{\tilde{A}}(x_i)/x_i$$

这里  $\mu_{\tilde{A}}(x_i)/x_i$  不是分数, 表示  $x_i$  对模糊集  $\tilde{A}$  的隶属程度是  $\mu_{\tilde{A}}(x_i)$ 。符号  $\sum$  不再表示数学求和, 而是各元素与其隶属度对应关系的一个总和。

(2) 设论域  $X$  是无限集, 则  $X$  上的一个模糊集  $\tilde{A}$  可以表示成

$$\tilde{A} = \int_{x \in X} \mu_{\tilde{A}}(x)/x$$

这里  $\int$  不再表示积分而是无穷逻辑和的意义。

### 1.3.2 基于目标函数的模糊聚类

在众多的模糊聚类算法中基于目标函数的模糊聚类由于其具有设计简单、解决问题的范围广, 最终可以归结为优化问题等优点而成为目前应用最广泛的模糊聚类算法。其中模糊 C-均值聚类是基于目标函数的聚类算法中最具代表性的聚类算法。

对于给定一个无标记的数据集  $X = \{x_1, x_2, \dots, x_n\} \subset \mathbb{R}^p$ , 将  $X$  分成  $c$  ( $1 < c < n$ ) 个模糊子集, 模糊 C-均值聚类处理以下优化问题:

$$\min_{(U,V)} J_m(U,V) = \min \text{Tr}(S_{fw}) \quad (1.1)$$

式(1.1)中,  $U \in M_{fc}$  是数据集  $X$  的模糊隶属度矩阵,  $U = [u_{ik}]_{c \times n}$ ;  $V = (v_1, v_2, \dots, v_c) \in \mathbb{R}^{cp}$ , 其中  $v_i \in \mathbb{R}^p$  是类中心矢量。 $M_{fc}$  是数据集  $X$  的模糊 C 划分空间:

$$M_{fc} = \left\{ U \in R^{cn} \mid u_{ik} = [0,1], \forall i,k; \sum_{i=1}^c u_{ik} = 1; 0 < \sum_{k=1}^n u_{ik} < n, \forall i \right\} \quad (1.2)$$

$S_{fw}$  为模糊类内散射矩阵(Fuzzy within Class Scatter Matrix):

$$S_{fw} = \sum_{i=1}^c \sum_{k=1}^n u_{ik}^m (x_k - v_i)(x_k - v_i)^T \quad (1.3)$$

则

$$\text{Tr}(S_{fw}) = \sum_{i=1}^c \sum_{k=1}^n u_{ik}^m D_{ik}^2 \quad (1.4)$$

这里  $D_{ik} = \|x_k - v_i\|$ ,  $c$  是要聚类的数目,  $n$  是聚类数据的数量,  $u_{ik}$  是数据  $x_k$  隶属于类  $i$  的隶属度值, 权重指数  $m \in (1, \infty)$ 。

为求解优化问题, 式(1.1)可用拉格朗日乘子法构造拉格朗日方程如下:

$$L(\alpha, u_{ik}, v_i) = \sum_{i=1}^c \sum_{k=1}^n u_{ik}^m \|x_k - v_i\|^2 - \alpha \left( \sum_{i=1}^c u_{ik} - 1 \right) \quad (1.5)$$

对方程(1.5)变量  $u_{ik}$  和  $v_i$  求偏导数后等于 0 得到:

$$\frac{\partial L(\alpha, u_{ik}, v_i)}{\partial u_{ik}} = mu_{ik}^{m-1} \|x_k - v_i\|^2 - \alpha = 0 \quad (1.6)$$

$$\frac{\partial L(\alpha, u_{ik}, v_i)}{\partial \alpha} = - \left( \sum_{i=1}^c u_{ik} - 1 \right) = 0 \quad (1.7)$$

$$\frac{\partial L(\alpha, u_{ik}, v_i)}{\partial v_i} = - 2 \sum_{k=1}^n u_{ik}^m (x_k - v_i) = 0 \quad (1.8)$$

求解以上式(1.6)~式(1.8)三个方程可得到如下方程:

$$u_{ik} = \left[ \sum_{j=1}^c \left( \frac{D_{ik}}{D_{jk}} \right)^{\frac{2}{m-1}} \right]^{-1}, \forall i, k \quad (1.9)$$

$$v_i = \frac{\sum_{k=1}^n u_{ik}^m x_k}{\sum_{k=1}^n u_{ik}^m}, \forall i \quad (1.10)$$

则模糊 C-均值聚类的迭代算法描述如下。

初始化:

(1) 固定  $m$  和  $c$  的值,  $n > c > 1, +\infty > m > 1$ ; 设置循环初始值  $r = 1$  和最大循

环次数为  $r_{\max}$ 。

(2) 选定初始类中心  $V^{(0)}$ :

循环:

步骤 1 用方程(1.9)更新隶属度值  $U^{(r)}$ ;

步骤 2 用方程(1.10)更新  $V^{(r)}$ ;

若满足条件( $\|V^{(r)} - V^{(r-1)}\| < \epsilon$ )或( $r > r_{\max}$ ), 则终止; 否则  $r = r + 1$ , 返回步骤 1。

### 1.3.3 模糊判别分析

模糊判别分析(FDA)是受 Fisher 线性判别分析的启发, 在模糊集基础上建立起来的有监督学习算法。在这里介绍 Chen 等提出模糊线性判别分析。

对于给定一个有标记的训练样本:

$$S = \{(x_1, y_1), \dots, (x_l, y_l)\} \subseteq (X \times Y)^l$$

这里  $X$  为输入空间,  $Y$  为输出空间,  $X \subseteq R^{l \times p}$ ,  $Y \subseteq R^p$ ,  $l$  为样本数, 样本  $x_i$  属于  $c$  个类别中的一个并且给定标记  $y_i \in \{1, 2, 3, \dots, c\}$ ,  $i=1, \dots, l$ 。

FDA 首先要将数据模糊化处理, 方法有:(1) 采用模糊 C-均值聚类算法。(2) 采用模糊支持向量机中的数据模糊化算法。(3) 采用模糊  $k$ -近邻法。FDA 通过求解以下广义瑞利商方程得到向量  $\omega$ :

$$\max J(\omega) = \frac{\omega^T S_{\text{BB}} \omega}{\omega^T S_{\text{TT}} \omega} \quad (1.11)$$

其中,  $S_{\text{BB}}$  为模糊类间散射矩阵(Fuzzy between Class Scatter Matrix);

$$S_{\text{BB}} = \sum_{i=1}^c \sum_{k=1}^l u_{ik}^m (v_i - \bar{x})(v_i - \bar{x})^T \quad (1.12)$$

$S_{\text{TT}}$  为模糊总体散射矩阵(Fuzzy Total Class Scatter Matrix):

$$S_{\text{TT}} = \sum_{i=1}^c \sum_{k=1}^l u_{ik}^m (x_k - \bar{x})(x_k - \bar{x})^T \quad (1.13)$$

式(1.13)和(1.12)中的  $\bar{x}$  是训练样本的平均值:

$$\bar{x} = \frac{1}{l} \sum_{j=1}^l x_j \quad (1.14)$$

求解方程(1.11)可以转化为求解下列特征方程:

$$S_{\text{TT}}^{-1} S_{\text{BB}} \beta = \lambda \beta \quad (1.15)$$

方程(1.15)的最大特征值就是方程(1.11)的解, 特征向量  $\beta$  对应特征值  $\lambda$ 。把测试样本  $z$  映射到向量  $\omega$  上:

$$w = z^T \omega \quad (1.16)$$

## 1.4 本章小结

本章主要介绍模糊分类的基础理论知识。介绍了模糊分类的发展概况，模糊聚类在国内外的研究概况。介绍了模糊分类的基本理论知识，包括模糊 C-均值聚类和模糊判别分析。

### 参 考 文 献

- [1] Zadeh L A. Fuzzy sets [J]. Information and Control, 1965, 8: 338 - 353.
- [2] Fisher R A. The use of multiple measurements in taxonomic problems [J]. Annals of Eugenics, 1936, 7(2): 179 - 188.
- [3] Foley D H, Sammon J W. An optimal set of discriminant vectors. IEEE Transaction on Computers, 1975, 24(3): 281 - 289.
- [4] Jin Z, Yang J Y, Hu Z S, Lou Z. Face recognition based on the uncorrelated discriminant transformation [J]. Pattern Recognition, 2001, 34: 1405 - 1416.
- [5] Bellman R, Kalaba R, Zadeh L. Abstraction and pattern classification [J]. Journal of Mathematical Analysis & Applications, 1996, 13(1): 1 - 7.
- [6] Wee W G. On generalizations of adaptive algorithms and application of the fuzzy sets concept to pattern classification [D]. Ph. D. thesis, Purdue University, Lafayette, Indiana, 1967.
- [7] Flake R H, Turner B L. Numerical classification for taxonomic problems [J]. Journal of Theoretical Biology, 1968, 20(2): 260 - 270.
- [8] Gitman I, Levine M. An algorithm for detecting unimodal fuzzy sets and its application as a clustering technique [J]. IEEE Transaction on Computers, 1970, C-19: 583 - 593.
- [9] Ruspini E. Numerical methods for fuzzy clustering [J]. Information Sciences, 1970 (6): 319 - 350.
- [10] Duda R O, Hart P E. Pattern Classification and Scene Analysis [M]. New York: Wiley, 1973.
- [11] 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.
- [12] Bezdek J C. Pattern recognition with fuzzy objective function algorithms [M]. New York: Plenum, 1981.
- [13] Bezdek J C, Pal M R, Keller J, Krisnapuram R. Fuzzy models and algorithms for pattern recognition and image processing [M]. Kluwer Academic, 1999.
- [14] Gustafson D E, Kessel W. Fuzzy clustering with a fuzzy covariance matrix [C]. In Proceedings of IEEE-CDC, 1979(2): 761 - 766.