

# 近邻分类方法及其应用

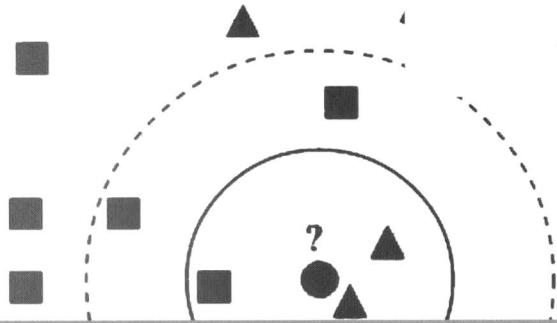
[下册]

郭躬德 陈黎飞 李南 ◎著

**Nearest Neighbour  
Classification Method  
and its Applications**



厦门大学出版社 国家一级出版社  
XIAMEN UNIVERSITY PRESS 全国百佳图书出版单位



# 近邻分类方法及其应用

[下册]

郭躬德 陈黎飞 李南 ◎著



厦门大学出版社 国家一级出版社  
XIAMEN UNIVERSITY PRESS 全国百佳图书出版单位

图书在版编目(CIP)数据

近邻分类方法及其应用·下册/郭躬德,陈黎飞,李南著.一厦门:厦门大学出版社,  
2014.4

ISBN 978-7-5615-5056-4

I. ①近… II. ①郭…②陈…③李… III. ①数据采掘-算法分析 IV. ①TP311.131

中国版本图书馆 CIP 数据核字(2014)第 068045 号



厦门大学出版社出版发行

(地址:厦门市软件园二期望海路 39 号 邮编:361008)

<http://www.xmupress.com>

xmup @ xmupress.com

厦门市明亮彩印有限公司印刷

2014 年 4 月第 1 版 2014 年 4 月第 1 次印刷

开本:787×1092 1/16 印张:15

插页:2 字数:341 千字

定价:38.00 元

本书如有印装质量问题请直接寄承印厂调换

# 前 言

计算机技术的普及应用给人类社会带来了深刻的变革,也使得我们所拥有的数据以前所未有的速度膨胀。随着大数据时代的到来,越来越多的人开始关注数据挖掘这一项大数据分析和处理的重要技术。作为数据挖掘的一种主要方法,分类(classification)——用于发掘隐藏在历史数据中的类别模式进而对未知事件做出预测或判断的技术——由于其很强的实用性成为了许多数据分析处理系统的基本构件。在机器学习领域,分类是有监督学习的代表性方法,其应用也已深入到信息检索、生物信息、客户关系管理等社会经济生活的方方面面。

近邻分类技术源来已久,可以追溯到早期对  $k_n$ -NN 规则(E. Fix 等人,1951 年)和  $k$ -NN 规则(T. Cover 等人,1967 年)的研究。随后发展起来的  $k$ NN 分类算法由于其原理简单、易于实现、可扩展性好、可解释性强等优点,备受青睐,位列 2005 年 ICDM 国际会议遴选的 10 大最有影响力的数据挖掘算法之一。如今,近邻分类方法无论在理论模型、算法还是在领域应用方面都吸引了众多的研究者和实践者,呈现出蓬勃发展的态势,涌现了一大批改进型新算法,也提出了一些基于近邻分类思想的新模型、新方法。“欲穷千里目,更上一层楼”,跟踪了解近年来取得的这些新进展才能进一步推进近邻分类技术的研究和深入应用,这也是本书出版的首要目的所在。

著者之一郭躬德教授系英国 Ulster 大学归国学者,长期从事近邻分类方法的研究与应用,早于 2003 年提出称为  $k$ NNModel 的近邻分类新方法,赢得业界热烈反响。近年来,在郭教授的带领下,福建师范大学数据挖掘与网络内容安全实验室开展了近邻分类理论方法与应用方面的系统研究,取得了一系列成果。在理论方法方面,研究团队提出了基于近邻思想的相似性度量新方法并将之推广到类属型数据,提出了增量学习、多代表点学习和

子空间近邻分类等新方法；应用研究涵盖了毒性物质预测、特征选择、文本分类以及数据流分类等近邻分类的新应用领域。本书将有关研究成果集结成册，以飨读者。

本书共六章，第一章介绍近邻分类算法及近年来的研究新进展，第二章至第六章每章介绍近邻分类的一类新方法。郭躬德主要编写本书的第二、三章，并参与编写第四、六章部分章节，约 20 万字；陈黎飞主要编写第四、五、六章，约 30 万字；李南主要编写第一章，并参与编写第四章部分章节，约 5 万字。研究生黄或、陈红、辛轶、黄杰、李南、张健飞、卢伟胜、兰天，访问学者陈雪云等参加了有关研究工作和部分章节写作。在写作过程中，参考了大量的国内外文献资料，在此一并表示感谢。

本书内容有误或不妥之处，欢迎读者批评指正。

郭躬德 陈黎飞 李南

福建师范大学数学与计算机科学学院

2014 年 1 月

# 目 录

## 上 册

第1章 近邻分类方法及其演变	1
1.1 分类概念、算法	1
1.2 经典的近邻分类方法及其演变	16
参考文献	24
第2章 近邻模型系列方法及其应用	29
2.1 $k$ 近邻模型分类算法	29
2.2 基于权重 $k$ 近邻模型的数据简化与分类	39
2.3 模糊 $k$ 近邻模型算法在可预测毒物学上的应用	50
2.4 最近邻分类的多代表点学习算法	62
2.5 改进的 $k$ 近邻模型方法在文本分类中的应用	72
2.6 部分模糊聚类的最近邻分类方法	87
参考文献	96
第3章 近邻模型的增量学习方法及其应用	102
3.1 基于 $kNN$ 模型的增量学习算法	102
3.2 增量 $kNN$ 模型的修剪策略研究	112
3.3 基于增量 $kNN$ 模型的分布式入侵检测架构	122
3.4 基于 $kNN$ 模型的层次纠错输出编码算法	131
参考文献	142

## 下 册

第 4 章 概念漂移数据流分类方法及其应用 .....	149
4.1 IKnnM-DHecoc:一种解决概念漂移问题的方法 .....	149
4.2 基于混合模型的数据流概念漂移检测 .....	164
4.3 面向高速数据流的集成分类器算法 .....	182
4.4 一种适应概念漂移数据流的分类算法 .....	192
4.5 基于少量类标签的概念漂移检测算法 .....	202
4.6 半监督层次纠错输出编码算法 .....	216
参考文献 .....	228
第 5 章 子空间近邻分类方法及其应用 .....	237
5.1 类依赖投影的文本分类方法 .....	237
5.2 多代表点的子空间分类算法 .....	253
5.3 基于投影原型的文本分类方法 .....	263
5.4 复杂数据的最优子空间分类方法 .....	280
5.5 基于特征子空间的概念漂移检测算法 .....	293
5.6 基于子空间集成的概念漂移数据流分类算法 .....	301
参考文献 .....	313
第 6 章 近邻方法的扩展及其应用 .....	320
6.1 基于空间覆盖的相似性度量及其对应的分类算法 .....	320
6.2 基于空间覆盖的相似性度量的特征选择算法 .....	333
6.3 基于空间覆盖的相似性度量的层次聚类算法 .....	340
6.4 基于类别子空间距离加权的互 $k$ 近邻算法 .....	347
6.5 针对类属性数据加权的 MKnn 算法 .....	356
6.6 属性加权的类属数据近邻分类 .....	364
参考文献 .....	378



















