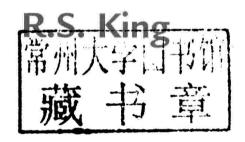
CLUSTER ANALYSIS AND DATA MINING

AN INTRODUCTION



CLUSTER ANALYSIS AND DATA MINING

An Introduction





MERCURY LEARNING AND INFORMATION

Dulles, Virginia

Boston, Massachusetts

New Delhi

Copyright ©2015 by Mercury Learning and Information LLC. All rights reserved.

This publication, portions of it, or any accompanying software may not be reproduced in any way, stored in a retrieval system of any type, or transmitted by any means, media, electronic display or mechanical display, including, but not limited to, photocopy, recording, Internet postings, or scanning, without prior permission in writing from the publisher.

Publisher: David Pallai
MERCURY LEARNING AND INFORMATION
22841 Quicksilver Drive
Dulles, VA 20166
info@merclearning.com
www.merclearning.com
1-800-758-3756

This book is printed on acid-free paper.

R.S. King. Cluster Analysis and Data Mining: An Introduction ISBN: 978-1-938549-38-0

The publisher recognizes and respects all marks used by companies, manufacturers, and developers as a means to distinguish their products. All brand names and product names mentioned in this book are trademarks or service marks of their respective companies. Any omission or misuse (of any kind) of service marks or trademarks, etc. is not an attempt to infringe on the property of others.

Library of Congress Control Number: 2014941165

141516321 Printed in the United States of America

Our titles are available for adoption, license, or bulk purchase by institutions, corporations, etc. Digital versions of this title are available at www.authorcloudware.com. Companion disc files may be obtained by contacting info@merclearning.com. For additional information, please contact the Customer Service Dept. at 1-800-758-3756 (toll free).

The sole obligation of MERCURY LEARNING AND INFORMATION to the purchaser is to replace the disc, based on defective materials or faulty workmanship, but not based on the operation or functionality of the product.

CLUSTER ANALYSIS AND DATA MINING

LICENSE, DISCLAIMER OF LIABILITY, AND LIMITED WARRANTY

By purchasing or using this book (the "Work"), you agree that this license grants permission to use the contents contained herein, but does not give you the right of ownership to any of the textual content in the book or ownership to any of the information or products contained in it. This license does not permit uploading of the Work onto the Internet or on a network (of any kind) without the written consent of the Publisher. Duplication or dissemination of any text, code, simulations, images, etc. contained herein is limited to and subject to licensing terms for the respective products, and permission must be obtained from the Publisher or the owner of the content, etc., in order to reproduce or network any portion of the textual material (in any media) that is contained in the Work.

MERCURY LEARNING AND INFORMATION ("MLI" or "the Publisher") and anyone involved in the creation, writing, or production of the companion disc, accompanying algorithms, code, or computer programs ("the software"), and any accompanying Web site or software of the Work, cannot and do not warrant the performance or results that might be obtained by using the contents of the Work. The author, developers, and the Publisher have used their best efforts to insure the accuracy and functionality of the textual material and/or programs contained in this package; we, however, make no warranty of any kind, express or implied, regarding the performance of these contents or programs. The Work is sold "as is" without warranty (except for defective materials used in manufacturing the book or due to faulty workmanship).

The author, developers, and the publisher of any accompanying content, and anyone involved in the composition, production, and manufacturing of this work will not be liable for damages of any kind arising out of the use of (or the inability to use) the algorithms, source code, computer programs, or textual material contained in this publication. This includes, but is not limited to, loss of revenue or profit, or other incidental, physical, or consequential damages arising out of the use of this Work.

The sole remedy in the event of a claim of any kind is expressly limited to replacement of the book, and only at the discretion of the Publisher. The use of "implied warranty" and certain "exclusions" vary from state to state, and might not apply to the purchaser of this product.

To the memory of my father, who made it possible and to LaJuan, the shining light in my life who survived the process.

PREFACE

his book is appropriate for a first course in clustering methods and data mining. Clustering and data mining methods are applicable in many fields of study, for example:

- 1. in the life sciences for developing complete taxonomies,
- 2. in the medical sciences for discovering more effective and economical means for making positive diagnosis in the treatment of patients,
- in the behavioral and social sciences for discerning human judgments and behavior patterns,
- 4. in the earth sciences for identifying and classifying geographical regions,
- in the engineering sciences for pattern recognition and artificial intelligence applications, and
- in decision and information sciences for analysis of markets and documents.

The first five chapters consider early historical clustering methods. Chapters 1 and 2 are an introduction to general concepts in clustering methods, with an emphasis on proximity measures and data mining. Classical numerical clustering methods are presented in Chapters 3 and 4: hierarchical and partitioned clustering. These methods are particularly defined

only on numeric data files. A clustering method implemented via multiple linear regression, judgmental analysis (JAN), is discussed in Chapter 5. JAN allows for numerical and categorical variables to be included in a clustering study.

All of the methods in Chapters 1 through 5 generate partitions on a study's data file, referred to as *crisp clustering* results. *Fuzzy clustering* methods presented in Chapter 6, capture partitions plus modified versions for the partitions. The modified partitions allow for overlapping clusters.

Chapter 7 is an introduction to the data mining topics of classification and association rules, which enable qualitative rather than simply quantitative data mining studies to be conducted.

Cluster analysis is essentially an art, but can be accomplished scientifically if the results of a clustering study can be validated. This is discussed in Chapter 8. Determination of the validity of individual clusters and the validation of a clustering, or collection of clusters, are discussed.

Chapter 9 surveys a variety of algorithms for clustering categorical data: ROCK, STIRR, CACTUS, and CLICK. These methods are dependent on underlying data structures and are applicable to relational databases.

Applications of clustering methods are presented in Chapters 10 through 11. Chapter ten discusses classical statistical methods for identifying outliers. Additionally, crisp and fuzzy clustering methods are applied to the outlier identification problem. Chapter 11 is an overview of model-based clustering. This is often used in physical science research studies for data generation.

A summary of the issues and trends in the cluster analysis field is made in Chapter 12. Besides giving recommendations for further study, an introduction to neural networks is presented. The appendices provide a variety of resources (software, URLs, algorithms, references) for the cluster analysis plus URLs for test data files.

The text is applicable to either a course on clustering, data mining, and classification or as a companion text for a first class in applied statistics. Clustering and data mining are good motivators and applications of the topics commonly included in an introductory applied statistics course.

The scheduling references for each of the chapters, in an applied statistics class, could be as follows:

Chapters 1-4: after study of descriptive statistics.

Chapter 9: immediately following Chapters 1-3.

Chapter 6: after study of descriptive statistics.

Chapter 10: after studying the Empirical Rule and Chebychev's Law.

Chapter 7: after studying probability.

Chapter 8: after study of hypothesis testing.

Chapter 5: after study of correlation, and both linear and multiple linear regression.

Chapter 11: after study of statistical inference.

No previous experience or background in clustering is assumed. Elementary statistics plus a brief exposure to data structures are the prerequisites. Informal algorithms for clustering data and interpreting results are emphasized. In order to evaluate the results of clustering and to explore data, graphical methods and data structures are used for representing data. Throughout the text, examples and references are provided, in order to enable the material to be comprehensible for a diverse audience.

CONTENTS

Pref	ace	value, value and			
Cha	ipte	er 1: Introduction to Cluster Analysis			
	1.1	What Is a Cluster?			
	1.2	Capturing the Clusters			
	1.3	Need for Visualizing Data			
	1.4	The Proximity Matrix			
	1.5	Dendrograms			
	1.6	Summary13			
	1.7	Exercises			
Chapter 2: Overview of Data Mining					
	2.1	What Is Data Mining?			
	2.2	Data Mining Relationship to Knowledge Discovery in Databases 19			
	2.3	The Data Mining Process			
	2.4	Databases and Data Warehousing			
	2.5	Exploratory Data Analysis and Visualization			
×	2.6	Data Mining Algorithms24			
	2.7	Modeling for Data Mining			

viii . Contents

2.8	Summary					
2.9	Exercises					
Chapter 3: Hierarchical Clustering						
3.1	$Introduction \ \dots \dots \dots$					
3.2	Single-Link versus Complet	e-Link Clustering				
3.3	Agglomerative versus Divisi	ve Clustering35				
3.4	Ward's Method					
3.5	Graphical Algorithms for Sic Complete-Link Clustering .	ngle-Link versus 				
3.6	Summary					
3.7	Exercises					
Chapter 4: Partition Clustering						
4.1	Introduction					
4.2	Iterative Partition Clusterin	g Method59				
4.3	The Initial Partition					
4.4	The Search for Poor Fits					
4.5	K-Means Algorithm					
	4.5.1 MacQueen's Method	d				
	4.5.2 Forgy's Method					
	4.5.3 Jancey's Method					
4.6	Grouping Criteria					
4.7	BIRCH, a Hybrid Method.					
4.8	Summary					
4.9	Exercises					
Chapter 5: Judgmental Analysis						
5.1	Introduction					
5.2	Judgmental Analysis Algorit	hm				
	5.2.1 Capturing R ²					
	5.2.2 Grouping to Optimi	ze Judges' R²				
	5.2.3 Alternative Method	for JAN89				
5.3	Judgmental Analysis in Res	earch				

5.4	Examp	le JAN Study	3
	5.4.1	Statement of Problem	3
	5.4.2	Predictor Variables96	6
	5.4.3	Criterion Variables	7
	5.4.4	Questions Asked	8
	5.4.5	Method Used for Organizing Data	8
	5.4.6	Subjects Judged103	3
	5.4.7	Judges	3
	5.4.8	Strategy Used for Obtaining Data103	3
	5.4.9	Checking the Model	6
	5.4.10	Extract the Equation	8
5.5	Summ	ary115	2
5.6	Exerci	ses	2
Chapt	er 6: Fu	uzzy Clustering Models and Applications 110	6
6.1	Introd	uction	6
6.2	The M	embership Function	1
6.3	Initial	Configuration	3
6.4	Mergin	ng of Clusters	4
6.5	Funda	mentals of Fuzzy Clustering	7
6.6	Fuzzy	C-Means Clustering	9
6.7	Induce	ed Fuzziness	7
6.8	Summ	ary14	1
6.9	Exerci	ses	2
Chapt	er 7: C	lassification and Association Rules 14	7
		uction	
7.2	Defini	ng Classification	8
		on Trees	
7.4	ID3 Ti	ree Construction Algorithm	2
	7.4.1	Choosing the "Best" Feature15	
¥	7.4.2	Information Gain Algorithm	5
	7.4.3	Tree Pruning	9

x • Contents

7.5 Bayesian Classification
7.6 Association Rules
7.7 Pruning
7.8 Extraction of Association Rules
7.9 Summary
7.10 Exercises
Chapter 8: Cluster Validity
8.1 Introduction
8.2 Statistical Tests
8.3 Monte Carlo Analysis
8.4 Indices of Cluster Validity
8.5 Summary
8.6 Exercises
Chapter 9: Clustering Categorical Data
9.1 Introduction
9.2 ROCK
9.3 STIRR
9.4 CACTUS241
9.5 CLICK
9.6 Summary
9.7 Exercises
Chapter 10: Mining Outliers
10.1 Introduction
10.2 Outlier Detection Methods
10.3 Statistical Approaches
10.4 Outlier Detection by Clustering
10.5 Fuzzy Clustering Outlier Detection
10.6 Summary
10.7 Exercises
Chapter 11: Model-based Clustering
11.1 Introduction
11.2 COBWEB: A Statistical and AI Approach

11.3 Mixture Model for Clustering
11.4 Farley and Raftery Gaussian Mixture Model286
11.5 Estimate the Number of Clusters
11.6 Summary
11.7 Exercises
Chapter 12: General Issues
12.1 Introduction
12.2 Data Cleansing
12.3 Which Proximity Measure Should Be Used?294
12.4 Identifying and Correcting Outliers294
12.5 Further Study Recommendations
12.6 Introduction to Neural Networks
12.7 Interpretation of the Results
12.8 Clustering "Correctness"?
12.9 Topical Research Exercises
On the DVD
Appendix A: Clustering Analysis with SPSS
Appendix B: Clustering Analysis with SAS
Appendix C: Neymann-Scott Cluster Generator Program Listing
Appendix D: Jancey's Clustering Program Listing
Appendix E: JAN Program
Appendix F: UCI Machine Learning Depository KD Nuggets Data Sets
Appendix G: Free Statistics Software (Calculator)
Appendix H: Solutions to Odd Exercises
Index

INTRODUCTION TO CLUSTER ANALYSIS

In This Chapter

- 1.1 What Is a Cluster?
- 1.2 Capturing the Clusters
- 1.3 Need for Visualizing Data
- 1.4 The Proximity Matrix
- 1.5 Dendrograms
- 1.6 Summary
- 1.7 Exercises

1.1 WHAT IS A CLUSTER?

Many of the decisions being made today involve more than one person. An important question in the group decision process is: "How does the group arrive at its final decision?" There have been a number of different mathematical and statistical approaches used by researchers attempting to model the decision-making process including game theory, information theory, and linear programming. Due to the large variety of decision-making situations, different types of decision processes, and the kinds of skills required, there is still a great deal of concern about the best way to make decisions. In many cases there is no objective approach. The individuals in the decision-making group each use their own set of criterion in reaching a decision. This approach might work in a situation where a consensus is

not needed. However, in the case where a single group decision is needed, there must be a "meeting of the minds."

One approach used is the *Delphi Technique*. This technique was designed in the early 1950s by the Rand Corporation to predict future outcomes. It is a group information gathering process to develop consensus opinion from a panel of experts on a topic of interest. In the normal Delphi scenario, the panel never meets face to face but interacts through questionnaires and feedback. This noncontact approach alleviates the worry over such issues as individual defensiveness or persuasiveness. However, opinions can be swayed due to a participant observing the responses of the rest of the panel. Another problem with the Delphi Technique is that the noncontact aspect is not feasible when, for example, the panel is the graduate admissions committee at a university.

Cluster analysis is another technique that has been used with success in the decision-making process. First, the investigator must determine the answer to "What is a cluster?" The **premise in cluster analysis** is: given a number of individuals, each of which is described by a set of numerical measures, devise a classification scheme for grouping the objects into a number of classes such that the objects within classes are *similar* in some respect and *unlike* those from other classes. These deduced classes are the clusters. The number of classes and the characteristics of each class must be determined from the data as discussed by Everett.¹

The key difference between cluster analysis and the Delphi Technique is that cluster analysis is strictly an objective technique. Whereas individual decisions can be swayed in an attempt to reach consensus in the Delphi process, or a "happy medium" is reached which does not really portray the feelings of the group as a whole. This is not the case in cluster analysis. Clusters of individuals are reached using an objective mathematical function. One particular type of cluster analysis called Judgmental ANalysis (JAN) takes the process one step further. Not only does it classify the panel into similar groups based on a related regression equation, but it also allows for these equations to be combined into a single policy equation. The JAN technique has been in use since the 1960s. It has proven to be an effective first step for methods of capturing and clustering the policies of judges.

Attempts at classification, that is sorting similar things into categories, can be traced back to primitive humans. The ability to classify is a necessary prerequisite for the development of language. Nouns, for example, are labels

¹ Everitt, B. S. (1980). Cluster analysis (2nd ed.). New York: Halsted Press.

used to classify a particular group of objects. Saying that a particular fourlegged animal is a "dog" allows us to put that animal into a category separate from cats, sheep, and horses. In other words, it allows us to communicate.

The classification of people and animals is almost as old as language. The early Hindus categorized humans into six types based on sex, physical, and behavioral characteristics. The early Greeks and Romans used classification to get a better understanding of the world around them. Galen, A.D. 129-199, defined nine temperamental types that were assumed to be related to a person's susceptibility to various diseases and to individual differences in behavior as discussed by Everitt. Development of a method to categorize animals into species was initiated by Aristotle. He started by dividing them into red blooded (vertebrates) and those not having red blood (invertebrates). He then subdivided the two groups again based on how their young were born. Theophrastus continued Aristotle's work, providing the groundwork for biological research for centuries. Eventually, new taxonomic systems were developed by such people as Linnaeus, Lindley, and Darwin. Classification was not restricted to the biological sciences. In chemistry, Mendeleyev used classification to develop the periodic tables, discussion by Everitt.

In the 1960s, two events led to an explosion of interest in cluster analysis. The availability and spread of large, high-speed computers opened up new possibilities for researchers. Additionally, the publication of *Principles of Numerical Taxonomy* by Sokal and Sneath² covered the following three important areas:

- 1. a number of different cluster analysis techniques
- 2. the use of computers in classification research
- 3. a radically empirical approach to biological taxonomy presented by Blashfield and Aldenderfer³

The need for cluster analysis arises in many fields of study. For example, Anderberg⁴ lists six areas where cluster analysis has been used successfully:

1. In the life sciences (biology, botany, zoology, etc.), the objects of analysis are life forms such as plants, animals, and insects. The

² Sokal, R. R., and Sneath, P. H. A. (1963). Principles of Numerical Taxonomy. W. H. Freeman.

³ Blashfield and Aldenderfer, M. S. (1978). The literature on cluster analysis. *Multivariate Behavioral Research*, 13, 271-295.

⁴ Anderberg, M. R. (1973). Cluster analysis for applications. New York: Academic Press.