Finding Groups in Data

An Introduction to Cluster Analysis

LEONARD KAUFMAN

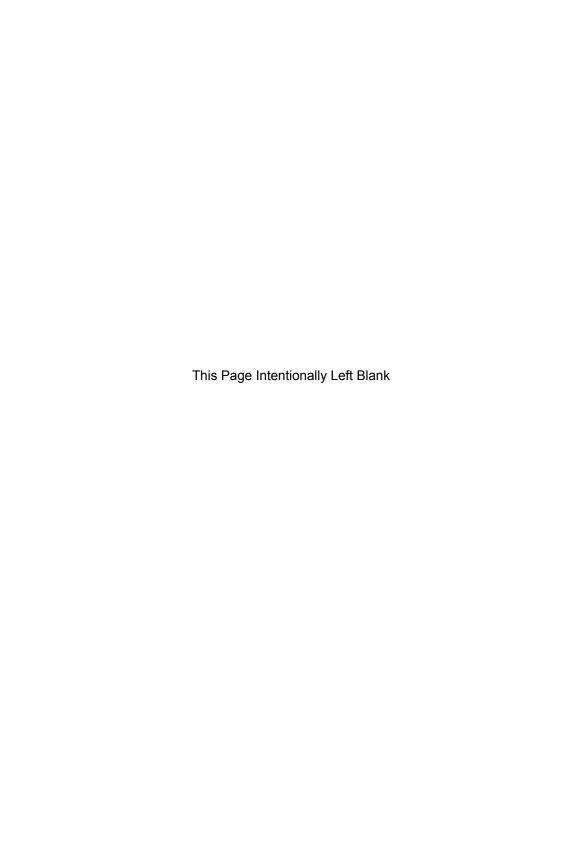
Vrije Universiteit Brussel, Brussels, Belgium

PETER J. ROUSSEEUW

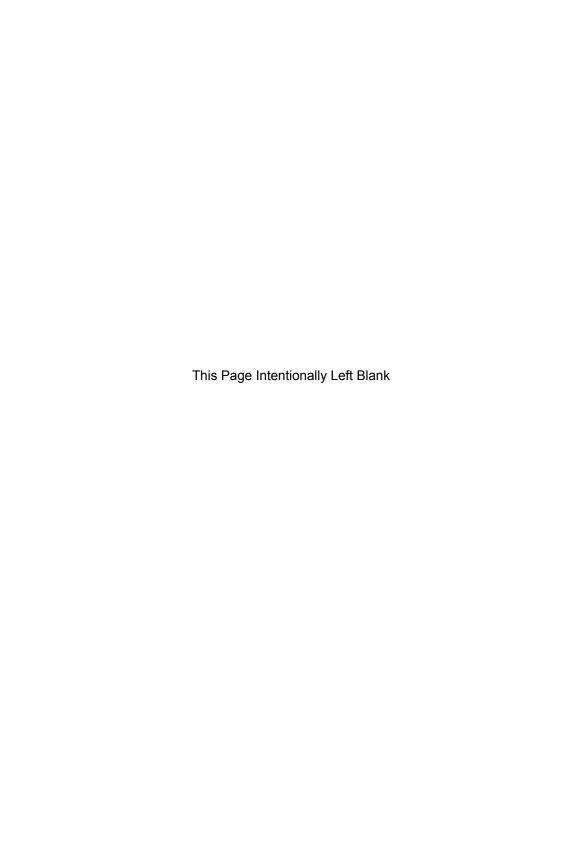
Universitaire Instelling Antwerpen, Antwerp, Belgium



A JOHN WILEY & SONS, INC., PUBLICATION



Finding Groups in Data



Finding Groups in Data

An Introduction to Cluster Analysis

LEONARD KAUFMAN

Vrije Universiteit Brussel, Brussels, Belgium

PETER J. ROUSSEEUW

Universitaire Instelling Antwerpen, Antwerp, Belgium



A JOHN WILEY & SONS, INC., PUBLICATION

Copyright © 1990, 2005 by John Wiley & Sons, Inc. All rights reserved.

Published by John Wiley & Sons, Inc., Hoboken, New Jersey. Published simultaneously in Canada.

No part of this publication may be reproduced, stored in a retrieval system or transmitted in any form or by any means, electronic, mechanical, photocopying, recording, scanning or otherwise, except as permitted under Sections 107 or 108 of the 1976 United States Copyright Act, without either the prior written permission of the Publisher, or authorization through payment of the appropriate per-copy fee to the Copyright Clearance Center, 222 Rosewood Drive, Danvers, MA 01923, (978) 750-8400, fax (978) 750-4470. Requests to the Publisher for permission should be addressed to the Permissions Department, John Wiley & Sons, Inc., 111 River Street, Hoboken, NJ 07030, (201) 748-6011, fax (201) 748-6008.

Limit of Liability/Disclaimer of Warranty: While the publisher and author have used their best efforts in preparing this book, they make no representation or warranties with respect to the accuracy or completeness of the contents of this book and specifically disclaim any implied warranties of merchantability or fitness for a particular purpose. No warranty may be created or extended by sales representatives or written sales materials. The advice and strategies contained herein may not be suitable for your situation. You should consult with a professional where appropriate. Neither the publisher nor author shall be liable for any loss of profit or any other commercial damages, including but not limited to special, incidental, consequential, or other damages.

For general information on our other products and services please contact our Customer Care Department within the U.S. at 877-762-2974, outside the U.S. at 317-572-3993 or fax 317-572-4002.

Wiley also publishes its books in a variety of electronic formats. Some content that appears in print, however, may not be available in electronic format.

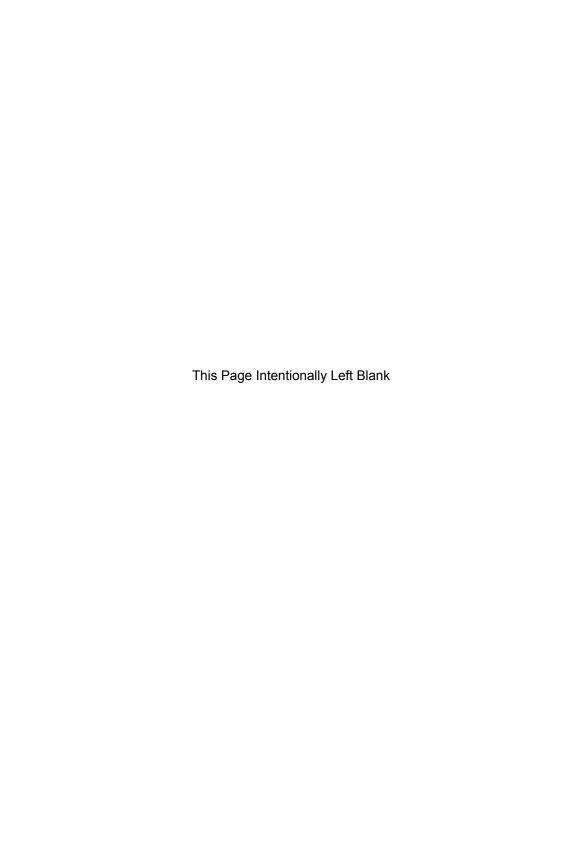
Library of Congress Cataloging-in-Publication is available.

ISBN 0-471-73578-7

Printed in the United States of America.

10 9 8 7 6 5 4 3 2 1

To Our Parents And To Jean Haezendonck



Preface

Cluster analysis is a very practical subject. Some 30 years ago, biologists and social scientists began to look for systematic ways to find groups in their data. Because computers were becoming available, the resulting algorithms could actually be implemented. Nowadays clustering methods are applied in many domains, including artificial intelligence and pattern recognition, chemometrics, ecology, economics, the geosciences, marketing, medical research, political science, psychometrics, and many more. This has led to a lot of different methods, and articles on clustering have appeared not only in statistical journals but also in periodicals of all these domains. Clustering is known under a variety of names, such as numerical taxonomy and automatic data classification.

Our purpose was to write an applied book for the general user. We wanted to make cluster analysis available to people who do not necessarily have a strong mathematical or statistical background. Rather than giving an extensive survey of clustering methods, leaving the user with a bewildering multitude of methods to choose from, we preferred to select a few methods that together can deal with most applications. This selection was based on a combination of methodological aims (mainly robustness, consistency, and general applicability) and our own experience in applying clustering to a variety of disciplines.

The book grew out of several courses on cluster analysis that we taught in Brussels, Delft, and Fribourg. It was extensively tested as a textbook with students of mathematics, biology, economics, and political science. It is one of the few books on cluster analysis containing exercises. The first chapter introduces the main approaches to clustering and provides guidance to the choice between the available methods. It also discusses various types of data (including interval-scaled and binary variables, as well as similarity data) and explains how these can be transformed prior to the actual

viii PREFACE

clustering. The other six chapters each deal with a specific clustering method. These chapters all have the same structure. The first sections give a short description of the clustering method, explain how to use it, and discuss a set of examples. These are followed by two sections (marked with * because they may be skipped without loss of understanding) on the algorithm and its implementation, and on some related methods in the literature. The chapters are relatively independent (except for Chapter 3 which builds on Chapter 2), allowing instructors to cover only one chapter in a statistics course. Another advantage is that researchers can pick out the method they need for their current application, without having to read other chapters. (To achieve this structure, some things had to be repeated in the text.) Occasionally, we handed a single chapter to someone working on a particular problem.

Chapters 2, 3, and 4 are about partitioning methods, whereas Chapters 5, 6, and 7 cover hierarchical techniques. Whenever possible, we constructed methods that cannot only analyze data consisting of measurements (i.e., objects with variables) but also data consisting of dissimilarities between objects. (This excluded parametric approaches, such as those based on multivariate mixture distributions.) We also wanted the methods to be consistent for large data sets. All the selected methods are of the L_1 type, which means that they minimize sums of dissimilarities (rather than sums of squared dissimilarities, as in the classical nonrobust methods). Some of the methods are new, such as the approach for partitioning large data sets and the L_1 method for fuzzy clustering. Also, the clusterings are accompanied by graphical displays (called silhouettes and banners) and corresponding quality coefficients, which help the user to select the number of clusters and to see whether the method has found groups that were actually present in the data.

Current statistical software packages contain only a few clustering techniques and they are not the more modern methods. This forced us to write new programs, the use of which is described in the book. Their present version is for IBM-PC and compatible machines, but the source codes are very portable and have run on several types of mainframes without problems. The programs (together with their sources and the data sets used in the book) are available on floppy disks by writing to the authors. The programs are also being integrated in the workstation package S-PLUS of Statistical Sciences, Inc., P.O. Box 85625, Seattle, WA 98145-1625.

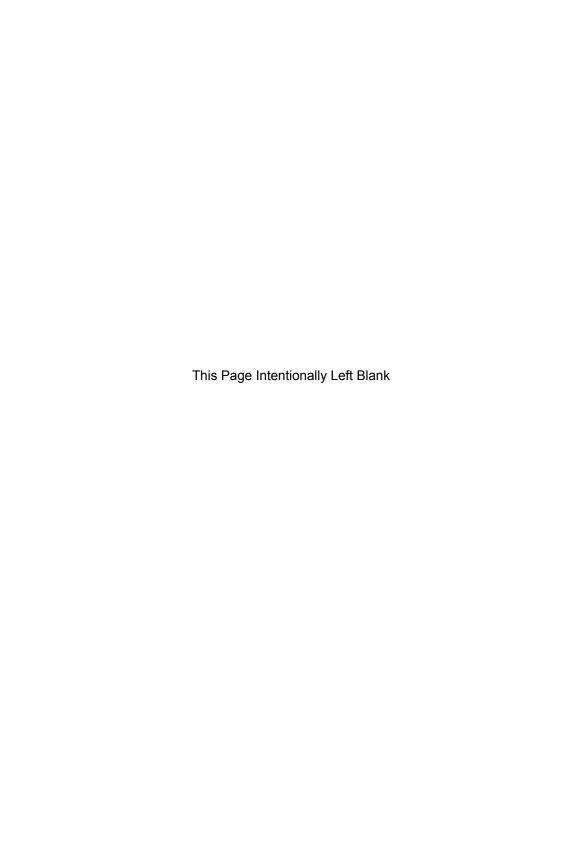
We are grateful to Frank Critchley, Jan de Leeuw, Gaetan Libert, Glenn Milligan, Frank Plastria, Marc Roubens, John Tukey, Bert van Zomeren, Howard Wainer, Michael Windham, and David Wishart for helpful suggestions and stimulating discussions on topics covered in this book and to

PREFACE ix

Etienne Trauwaert for contributing to Section 4 of Chapter 4. Valuable comments on the manuscript were given by Phil Hopke, Doug Martin, Marie-Paule Derde, and two reviewers. Annie De Schrijver was responsible for drawing several figures. Finally, we would like to thank our wives, Jacqueline and Lieve, for their patience and support.

LEONARD KAUFMAN PETER J. ROUSSEEUW

Edegem, Belgium August, 1989



Contents

1. Introduction

	1.	Motivation, 1	
	2.	Types of Data and How to Handle Them, 3	
		2.1 Interval-Scaled Variables, 4	
		2.2 Dissimilarities, 16	
		2.3 Similarities, 20	
		2.4 Binary Variables, 22	
		2.5 Nominal, Ordinal, and Ratio Variables, 28	
		2.6 Mixed Variables, 32	
	3.	Which Clustering Algorithm to Choose, 37	
		3.1 Partitioning Methods, 38	
		3.2 Hierarchical Methods, 44	
	4.	A Schematic Overview of Our Programs, 50	
	5.	Computing Dissimilarities with the Program DAISY, 52	
	Exe	rcises and Problems, 63	
2.	Part	titioning Around Medoids (Program PAM)	68
	1.	Short Description of the Method, 68	
	2.	How to Use the Program PAM, 72	
		2.1 Interactive Use and Input, 72	
		2.2 Output, 80	
		2.3 Missing Values, 88	
	3.	Examples, 92	
		More on the Algorithm and the Program, 102	
		4.1 Description of the Algorithm, 102	
		4.2 Structure of the Program, 104	
		•	•
			хi

1

xii CONTENTS

3.

4.

* 5.	Related Methods and References, 108				
	5.1 The k-Medoid Method and Optimal Plant Location, 10	8			
	5.2 Other Methods Based on the Selection of Representa				
	Objects, 110				
	5.3 Methods Based on the Construction of Central Points,	111			
	5.4 Some Other Nonhierarchical Methods, 116				
	5.5 Why Did We Choose the k-Medoid Method?, 117				
	5.6 Graphical Displays, 119				
Exe	rcises and Problems, 123				
Clu	stering Large Applications (Program CLARA)	126			
1.	Short Description of the Method, 126				
2.	•				
	2.1 Interactive Use and Input, 127				
	2.2 Output, 130				
	2.3 Missing Values, 134				
3.	An Example, 139				
*4.	More on the Algorithm and the Program, 144				
	4.1 Description of the Algorithm, 144				
	4.2 Structure of the Program, 146				
	4.3 Limitations and Special Messages, 151				
	4.4 Modifications and Extensions of CLARA, 153				
*5.	Related Methods and References, 155				
	5.1 Partitioning Methods for Large Data Sets, 155				
	5.2 Hierarchical Methods for Large Data Sets, 157				
	5.3 Implementing CLARA on a Parallel Computer, 160				
Exe	ercises and Problems, 162				
Fuz	zy Analysis (Program FANNY)	164			
1.	The Purpose of Fuzzy Clustering, 164				
2.	How to Use the Program FANNY, 166				
	2.1 Interactive Use and Input, 167				
	2.2 Output, 170				
3.					
*4.	More on the Algorithm and the Program, 182				
	4.1 Description of the Algorithm, 182				
	4.2 Structure of the Program, 188				

CONTENTS xiii

	⁻ ⊃.	Related Methods and References, 189	
		5.1 Fuzzy k-Means and the MND2 Method, 189	
		5.2 Why Did We Choose FANNY?, 191	
		5.3 Measuring the Amount of Fuzziness, 191	
		5.4 A Graphical Display of Fuzzy Memberships, 195	
	Exe	rcises and Problems, 197	
5.	Agglomerative Nesting (Program AGNES)		
	1.	Short Description of the Method, 199	
	2.	How to Use the Program AGNES, 208	
		2.1 Interactive Use and Input, 2082.2 Output, 209	
	3.	Examples, 214	
	*4.	More on the Algorithm and the Program, 221	
		4.1 Description of the Algorithm, 221	
		4.2 Structure of the Program, 223	
	* 5.	Related Methods and References, 224	
		5.1 Other Agglomerative Clustering Methods, 224	
		5.2 Comparing Their Properties, 238	
		5.3 Graphical Displays, 243	
	Exe	rcises and Problems, 250	
6.	Divisive Analysis (Program DIANA)		253
	1.	Short Description of the Method, 253	
	2.	How to Use the Program DIANA, 259	
	3.	Examples, 263	
	*4.	More on the Algorithm and the Program, 271	
		4.1 Description of the Algorithm, 271	
		4.2 Structure of the Program, 272	
	* 5.	Related Methods and References, 273	
		5.1 Variants of the Selected Method, 273	
		5.2 Other Divisive Techniques, 275	
	Exe	rcises and Problems, 277	
7.	Monothetic Analysis (Program MONA)		
	1.	Short Description of the Method, 280	
	2.	How to Use the Program MONA, 283	
		5 - 1	