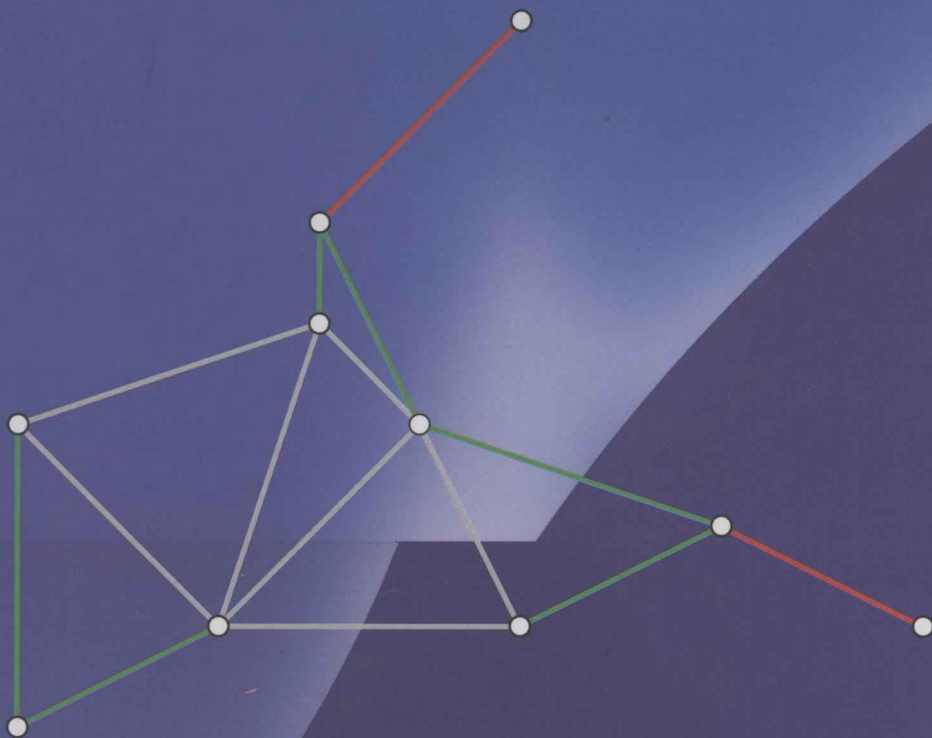


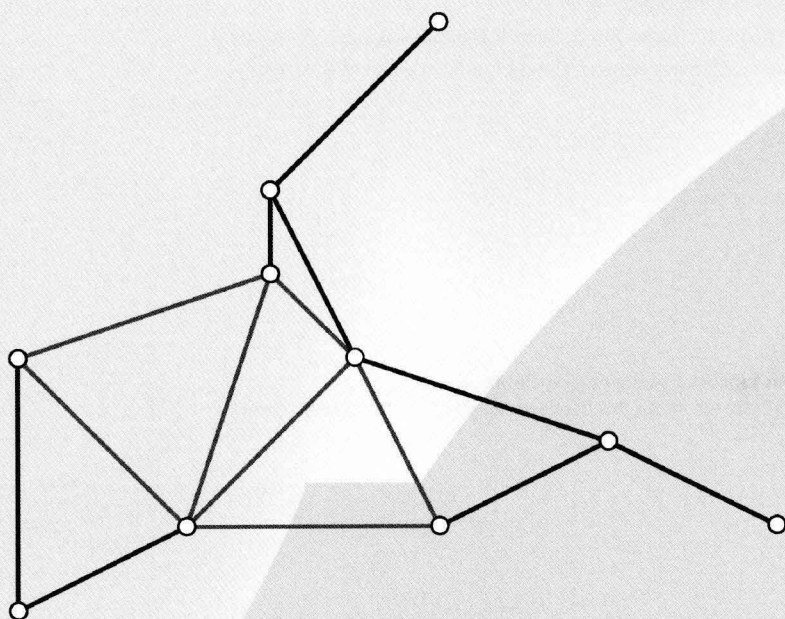
Interdisciplinary Mathematical Sciences – Vol. 10



# Ordinal and Relational Clustering

Melvin F Janowitz

**Interdisciplinary Mathematical Sciences – Vol. 10**



# **Ordinal and Relational Clustering**

**Melvin F Janowitz**

Rutgers University, USA

 **World Scientific**

NEW JERSEY • LONDON • SINGAPORE • BEIJING • SHANGHAI • HONG KONG • TAIPEI • CHENNAI

*Published by*

World Scientific Publishing Co. Pte. Ltd.

5 Toh Tuck Link, Singapore 596224

*USA office:* 27 Warren Street, Suite 401-402, Hackensack, NJ 07601

*UK office:* 57 Shelton Street, Covent Garden, London WC2H 9HE

**British Library Cataloguing-in-Publication Data**

A catalogue record for this book is available from the British Library.

**Interdisciplinary Mathematical Sciences — Vol. 10**

**ORDINAL AND RELATIONAL CLUSTERING**

**(With CD-ROM)**

Copyright © 2010 by World Scientific Publishing Co. Pte. Ltd.

*All rights reserved. This book, or parts thereof, may not be reproduced in any form or by any means, electronic or mechanical, including photocopying, recording or any information storage and retrieval system now known or to be invented, without written permission from the Publisher.*

For photocopying of material in this volume, please pay a copying fee through the Copyright Clearance Center, Inc., 222 Rosewood Drive, Danvers, MA 01923, USA. In this case permission to photocopy is not required from the publisher.

ISBN-13 978-981-4287-20-3

ISBN-10 981-4287-20-2

Printed in Singapore by Mainland Press Pte Ltd.

# **Ordinal and Relational Clustering**

## INTERDISCIPLINARY MATHEMATICAL SCIENCES

**Series Editor:** Jinqiao Duan (*Illinois Inst. of Tech., USA*)

**Editorial Board:** Ludwig Arnold, Roberto Camassa, Peter Constantin,  
Charles Doering, Paul Fischer, Andrei V. Fursikov, Xiaofan Li,  
Sergey V. Lototsky, Fred R. McMorris, Daniel Schertzer,  
Bjorn Schmalfuss, Yuefei Wang, Xiangdong Ye, and Jerzy Zabczyk

---

### *Published*

- Vol. 1: Global Attractors of Nonautonomous Dissipative Dynamical Systems  
*David N. Cheban*
- Vol. 2: Stochastic Differential Equations: Theory and Applications  
A Volume in Honor of Professor Boris L. Rozovskii  
*eds. Peter H. Baxendale & Sergey V. Lototsky*
- Vol. 3: Amplitude Equations for Stochastic Partial Differential Equations  
*Dirk Blömker*
- Vol. 4: Mathematical Theory of Adaptive Control  
*Vladimir G. Sragovich*
- Vol. 5: The Hilbert–Huang Transform and Its Applications  
*Norden E. Huang & Samuel S. P. Shen*
- Vol. 6: Meshfree Approximation Methods with MATLAB  
*Gregory E. Fasshauer*
- Vol. 7: Variational Methods for Strongly Indefinite Problems  
*Yanheng Ding*
- Vol. 8: Recent Development in Stochastic Dynamics and Stochastic Analysis  
*eds. Jinqiao Duan, Shunlong Luo & Caishi Wang*
- Vol. 9: Perspectives in Mathematical Sciences  
*eds. Yisong Yang, Xinchu Fu & Jinqiao Duan*
- Vol. 10: Ordinal and Relational Clustering (with CD-ROM)  
*Melvin F. Janowitz*

For Trudy

Who makes all things possible!

# Preface

***Underlying motivation for the book:*** *If data has only ordinal significance, then taking averages is not a meaningful basis for comparison.* This is not a new observation. There is a deep theory of measurement, and a substantial literature involving this and related issues. The interested reader might consult sources like Refs. [46, 55, 56].

This book is rather unusual in its organization as well as in its purpose. It is neither a text book nor is it a well documented exhaustive research volume. A middle ground is sought. Its subject is a type of exploratory data analysis that has come to be known as *cluster analysis*. Despite claims to the contrary in the literature, arguments will be presented to show that the input data to a clustering algorithm might have nothing more than some version of ordinal significance. For that reason, a careful approach will be presented for the analysis of ordinal data. Though mathematicians often try to develop general properties of whatever system they are studying, it seems clear that there is no axiom system that will provide a “best” method of discovering any internal structure of an arbitrary data set. For that reason, we will present properties that cluster algorithms may or may not enjoy, rather than axioms that should necessarily hold for every algorithm and all data.

Where appropriate, we shall illustrate abstract ideas with concrete numerical examples, as well as with careful mathematical arguments. Much of our treatment is based on ideas originally appearing in Ref. [44]. For readers wishing a more comprehensive introduction to cluster analysis, possible references include Refs. [28, 30, 32, 51, 63]. Though the treatment is mathematical in nature, it will be written so as to make it accessible to researchers interested in concrete data analysis. The first five chapters deal with standard cluster analysis, while the remaining two chapters present a new and more general theory that subsumes both cluster analysis and also other related techniques of data analysis. In particular, it includes some aspects of formal concept analysis (Refs. [18, 26]), and makes contact with symbolic data analysis [13].

A *clustering problem* may take as its input a finite data set  $P$  having at least three elements, together with a list of finitely many attributes that the data might have in varying degrees. Alternately, it may consist of a numerical (perhaps ordinal)

measure of similarity or dissimilarity between pairs of objects of  $P$ . Though the ultimate output is to discover some sort of internal classification structure that  $P$  may have, we argue that an appropriate intermediate step is to form a dissimilarity coefficient that somehow summarizes this internal structure. The key idea is that the sought after output structure must be determined entirely by the input, and is in no way directly dependent on any external criteria. We do not argue that external criteria should be ignored, but only take the view that when such considerations are made, we are no longer in the realm of cluster analysis. We shall largely ignore the passage from attribute data to similarity or dissimilarity measures, and concentrate instead on the transformation of a dissimilarity measure into some sort of classification or nested sequence of classifications by means of first creating a “summary” measure of dissimilarity from which the ultimate output may be constructed. Thus, we usually take our input for a cluster algorithm to be a dissimilarity coefficient. The reader should note that the terms “dissimilarity coefficient” and “dissimilarity measure” are being used interchangeably. The abbreviation DC provides a convenient way to refer to them.

Research in the social and behavioral sciences is replete with examples where perceptions or preferences are first quantified, and then on the basis of this quantification, a classification or a decision must be reached. For example, judges will give numerical ratings to paintings or to athletic performances. Though these ratings are subjective with criteria that vary from judge to judge, one hopes that if the rating for A is less than the rating for B, then in some interpretable sense, B is superior to A. In other words, even though the numbers themselves may not have meaning, the fact that the rating for A is less than the rating for B may still have significance. Much the same situation applies in general to cluster analysis. We will discuss algorithms whose input is a numerical measure of dissimilarity. If one knew the nature of the data and just how the measure  $d$  of dissimilarity was constructed, the actual values of  $d(x, y)$  might be interpretable. But we cannot assume such knowledge. For a measure  $d$  of unspecified origin, the most one can hope for, and all that is claimed, is the fact that  $d(x, y) < d(s, t)$  should imply that  $x$  is more similar to  $y$ , than  $s$  is to  $t$ . We will briefly consider the behavior of cluster algorithms when faced with noisy data. Here versions of continuity are a natural consideration, as are the nature of various properties of the input data implied by the choice of cluster algorithm. Naturally, there will be other considerations.

Here is an intuitive road map that will guide your journey through the text. We begin by introducing the underlying ideas by means of examples and informal discussions. This is followed by a more mathematical treatment of the same material, again illustrated by examples. The general theme we follow is that of looking at the action of various mappings on the reals with a numerical dissimilarity measure  $d$ . Thus we carefully study the effects of the transformation of the dissimilarity coefficient  $d$  into the dissimilarity coefficient  $\theta d$ , where  $\theta$  is some mapping on the nonnegative reals. This is the content of the third chapter, while Chapter 4



considers ideas related to continuity. It turns out to be possible to classify cluster methods according to how they interact with various mappings on the reals. The idea is that this classification tells us something about the data assumptions that are implicitly made by the choice of cluster method. Such is the content of the fifth chapter.

Chapter 6 begins with an introduction to formal concept analysis, thereby laying a foundation for making a connection between cluster analysis and this seemingly unrelated discipline. The reader is shown situations where dissimilarities are measured in settings more general than the real number system. Dissimilarities taking values in a partially ordered set  $L$  will be of special interest. It is argued here that it is appropriate to actually view the dissimilarity between  $x$  and  $y$  as an *order filter* of the poset  $L$  in which dissimilarities are measured. This leads to a dramatically different view of a dissimilarity coefficient  $d$ . For objects  $x, y$ , the idea is that  $d(x, y)$  should denote the set of levels of the poset  $L$  at which  $\{x, y\}$  is a *candidate* for clustering. A cluster method  $F$  decides at each level which cluster candidates should actually be merged to form clusters. The output of  $F$  is then a dissimilarity coefficient whose image takes values in the order filters of  $L$ . Associated with this output is a collection of partitions of subsets of the underlying set, thus making a natural connection with the type of display used in formal concept analysis. This is a display of some or all of the possible clusters that arise. The goal of this type of clustering is to provide possible hypotheses for the internal structure of  $P$ . Thus rather than assume that there is a unique “true” hierarchical structure that we are trying to estimate, we make no statistical assumptions, and instead take the view that we are seeking possible hidden internal structures of the given data set. Note that this approach allows us to consider a DC taking values in the attribute space associated with the input data. For that reason, we will call this new general discipline “relational clustering”. Naturally, we do not preclude the possibility of an underlying “true” structure that we are estimating, but our goal is to present a model that is also suitable for use with data mining applications. The actual connection with FCA is made by considering dissimilarities taking values in a finite Boolean algebra.

The final chapter presents a more formal and more general view of this material, and illustrates the idea by presenting the results of some clustering algorithm based on these ideas. In particular, this allows one to consider attributes whose values vary within the confines of the individual objects being clustered. This enables contacts with percentile clustering, confidence interval clustering, and symbolic data analysis. It is hoped that this introduction to such a general model of cluster analysis will inspire further development of this new subject.

We close with a word about organization and terminology. The text is organized into chapters and sections. Numbering will be by divisions, so that a reference of the form Theorem 2.3.5 would indicate Theorem 5 of Section 3 of Chapter 2. Though we try to make the treatment as self-contained as possible, we do assume a basic

knowledge of set theory, functions, and the real number system. Sections prefaced with a  $\star$  should be viewed as optional and may safely be omitted. These sections tend to be technical and require a certain amount of extra sophistication on the part of the reader. They may also contain ideas that are not often of interest to the typical user of clustering techniques. The terminology and notation are (at least to the author) standard. There is included a list of tables, a list of figures, a list of symbols, as well as a detailed index for the text. You will find a CD containing MATLAB software accompanying the book. This software is designed to allow the reader to personally check the illustrative examples in the book; it will be discussed in Sections 3.4.3, 6.5, 7.4, 7.5 and 7.6. Any needed corrections or improvements will be placed in the website

<http://www.worldscibooks.com/mathematics/7449.html>

## Acknowledgments

Before proceeding, the author expresses his gratitude to CLUSTAN Ltd. for granting permission to use the ClustanGraphics clustering package in connection with processing some of the illustrative examples. Illustrative software is provided that makes use of version 7.9 of MATLAB<sup>®</sup> and the Optimization Toolbox. Product information may be obtained from

The MathWorks, Inc.  
3 Apple Hill Drive  
Natick, MA 01760-2098 USA  
Tel: 508-647-7000, Fax: 508-647-7001  
E-mail: [info@mathworks.com](mailto:info@mathworks.com)  
Web: [www.mathworks.com](http://www.mathworks.com)

We especially thank Ed Scheinerman for making available the Matgraph Toolbox <http://www.ams.jhu.edu/~ers/matgraph> that was used as a basis for many of our MATLAB<sup>®</sup> calculations.

Professor Thomas Riedel, Mathematics Department, University of Louisville taught a graduate course based on an early version of the text during the Summer of 2008. Both he and his class made valuable suggestions that improved the quality of the exposition. We especially thank his students for their input. Among them are Richard Lagani, Amanda Ledford, and Christy Parks. We also thank Professors Robert C. Powers (University of Louisville), Richard Greechie (Louisiana Polytechnic University), and Li Chen (University of the District of Columbia) for their comments and suggestions. Finally, we thank Ralph Freese (University of Hawaii) for modifying and making available for use with the book his software package *LatDraw* for drawing lattices, (<http://www.latdraw.org>). We also thank the Computer Science Laboratory at the University of the District of Columbia for a version of the software that was used to produce Figures 7.1, 7.2, and 7.3 of the text. A big thanks goes to Fred McMorris (Illinois Institute of Technology) both for his editorial role with World Scientific, and his encouragement over many years as a valued research colleague. We also owe a debt of gratitude to the editor of the series in which this volume appears: Professor Jinqiao (Jeffrey) Duan (Illinois Institute of Technology), as well as to the editorial staff of World Scientific: notably to Rajesh Babu and Rok Ting Tan.

# Contents

<i>Preface</i>	vii
<i>Acknowledgments</i>	xi
<i>List of Figures</i>	xv
<i>List of Tables</i>	xvii
1. Informal background	1
1.1 Types of relations . . . . .	2
1.2 Nine integer example . . . . .	5
1.3 Mammal milk example . . . . .	7
2. Dissimilarities and clusters	13
2.1 Getting more formal . . . . .	14
2.2 Dendrograms and ultrametrics . . . . .	21
2.3 The Jardine-Sibson model . . . . .	24
2.4 Worked examples . . . . .	29
2.5 ★ Characterization of the family of clusters . . . . .	39
2.6 ★ ML-sets and complete-linkage clustering . . . . .	45
3. Ordinal data	49
3.1 Monotone equivariance . . . . .	49
3.2 Actions of isotone mappings . . . . .	53
3.3 Flat and 0-flat cluster methods . . . . .	58
3.4 Clustering based on bridges . . . . .	62
3.4.1 Bridges and generalized bridges . . . . .	63
3.4.2 Bridge removal clustering . . . . .	65
3.4.3 Software considerations . . . . .	68

4.	Continuity and ordinal continuity	75
4.1	Notions of Limits . . . . .	75
4.2	Continuity . . . . .	79
4.3	★ Order continuity . . . . .	86
5.	Classification of monotone equivariant cluster methods	93
5.1	Classification kernels . . . . .	93
5.2	The clustering connection . . . . .	103
5.3	Kernels of residuated mappings . . . . .	112
5.3.1	Residuated and residual mappings . . . . .	112
5.3.2	Kernels for residuated mapping on $\mathbb{R}_0^+$ . . . . .	117
5.4	Normalized dissimilarities . . . . .	119
5.5	Clustering connection for normalized residuated mappings . . . . .	122
6.	Clustering based on posets	129
6.1	Galois connections . . . . .	129
6.2	Formal concept analysis . . . . .	130
6.3	Boolean dissimilarities . . . . .	136
6.4	Relational clustering . . . . .	143
7.	A new poset model	145
7.1	An informal approach . . . . .	145
7.2	A more formal approach . . . . .	146
7.3	Algorithms in the poset model . . . . .	150
7.4	★ The software . . . . .	153
7.5	* The Majority induced order . . . . .	158
7.6	Confidence interval clustering . . . . .	162
7.7	Thoughts for the future . . . . .	167
	<i>List of Symbols</i>	169
	<i>Bibliography</i>	171
	<i>Index</i>	175

# List of Figures

2.1	Threshold relations for Example 2.5 . . . . .	17
2.2	Clusters for Example 2.5 . . . . .	18
2.3	Fundamental depiction of a dissimilarity coefficient . . . . .	18
2.4	Step function display of an NSC . . . . .	19
2.5	Two visualizations of a dendrogram . . . . .	24
2.6	Recapturing $Tu(h)$ . . . . .	24
2.7	Single-linkage for the nine integers . . . . .	40
2.8	Complete-linkage clustering for the nine integers . . . . .	40
2.9	Average-linkage clustering for the nine integers . . . . .	40
2.10	Jardine-Sibson complete-linkage clustering for the nine integers . . . . .	40
2.11	Average-linkage clustering for reversed order of entry . . . . .	40
2.12	Proportional-linkage clustering with $u = 0.7$ . . . . .	40
2.13	Illustration of maximal cliques . . . . .	46
3.1	Action of $\theta_d$ and $\theta_d^+$ . . . . .	57
3.2	Bridges and generalized bridges . . . . .	65
3.3	Single-linkage for the sixteen integers . . . . .	67
3.4	Bridge clustering illustration . . . . .	67
3.5	Generalized bridge deletions at level 5 . . . . .	67
3.6	Generalized bridge deletions at level 8 . . . . .	67
4.1	Limit of a sequence . . . . .	77
4.2	Sequence with no limit . . . . .	77
4.3	Limit from the right of the sequence $(\theta_n)$ is $\tau_0$ . . . . .	78
4.4	Limit from the left of $(\psi_n)$ is $\tau_0$ . . . . .	78
4.5	Limit from the right of $(\mu_n)$ is $\tau_k$ ( $k = 1.25$ ) . . . . .	78
4.6	Limit from the left of $(\zeta_n)$ is $\nu_j$ ( $j = 1.75$ ) . . . . .	78
4.7	Single-linkage clustering . . . . .	81
4.8	(a) Complete-linkage clustering (b) Average-linkage clustering . . . . .	82

5.1	The lattice of dependence functions for 0-preserving isotone mappings on $\mathfrak{R}_0^+$ . . . . .	102
5.2	The lattice of compatability classes for ME cluster methods . . . . .	107
5.3	The lattice of dependence functions for residuated mappings on $\mathfrak{R}_0^+$ . . .	119
5.4	The lattice of kernels for $\mathfrak{R}_1$ . . . . .	123
5.5	The lattice of normalized ME functions for $\mathfrak{R}_1$ . . . . .	126
6.1	Nonempty extents of the nine integer example . . . . .	134
6.2	Nonempty extents of the mammal milk example . . . . .	135
6.3	Clusters for the nine integer example with $d_2(x, x) = 0$ . . . . .	140
7.1	General complete-linkage clustering on crime data . . . . .	155
7.2	General single-linkage clustering on milk data . . . . .	159
7.3	General complete-linkage on the mammal milk data . . . . .	159
7.4	Confidence Intervals . . . . .	163
7.5	Single-linkage clustering on average Vietnam combat deaths . . . . .	166
7.6	Complete-linkage clustering on average Vietnam combat deaths . . . . .	166

# List of Tables

1.1	Examples of relations on $X = \{a, b, c\}$ . . . . .	3
1.2	Names of attributes . . . . .	5
1.3	Attributes for the nine integer example . . . . .	5
1.4	Simple matching coefficient for the nine integer example . . . . .	6
1.5	Jaccard coefficient for the nine integer data . . . . .	6
1.6	Percentage composition of mammal milk . . . . .	9
1.7	Chebychef DC for the mammal milk data . . . . .	9
1.8	Manhattan distance for the mammal milk data . . . . .	10
1.9	Squared Euclidean distance for the mammal milk data . . . . .	10
1.10	Z-score version of the mammal milk data . . . . .	10
1.11	Mammal milk data standardized to have range $[0,1]$ . . . . .	11
1.12	Mammal milk data standardized to have mean 1 . . . . .	11
2.1	Illustration of simple matching coefficient for the nine integer example .	35
2.2	Single-linkage output for Table 2.1 . . . . .	36
2.3	Complete-linkage output for Table 2.1 . . . . .	37
2.4	Average-linkage output for Table 2.1 . . . . .	39
2.5	Proportional-linkage output with $u = 0.7$ for Table 2.1 . . . . .	39
3.1	16 points in the Euclidean plane . . . . .	66
5.1	Types of $q$ -intervals . . . . .	97
5.2	The dependence functions for the seven classes of $Q_{\mathfrak{M}}$ for $\theta \in \mathfrak{M}$ . . . .	102
5.3	The Classes of ME cluster methods . . . . .	103
5.4	The seven classes of $K(\theta)$ for $\theta \in \mathfrak{M}$ . . . . .	118
5.5	The five classes of $K(\theta)$ for $\theta$ residuated on $\mathfrak{R}_0^+$ . . . . .	118
5.6	Types of $q$ -intervals . . . . .	120
5.7	The nine classes for 1-reserving residuated mappings on $[0,1]$ . . . . .	123
5.8	Dependence functions for the eleven classes of normalized monotone equivariant cluster functions . . . . .	127



6.1	The formal concepts of $(\mathfrak{G}, \mathfrak{M}, \perp)$ . . . . .	132
6.2	Properties of the nine integers . . . . .	132
6.3	Attributes of the first nine integers . . . . .	133
6.4	Percentage composition of mammal milk . . . . .	134
6.5	Binary version of mammal milk example . . . . .	135
6.6	Illustration of $d_1$ and $d_2$ . . . . .	139
6.7	Illustration of the $d_1$ coefficient for the nine integers . . . . .	140
6.8	$d_1$ clusters at level 11100 . . . . .	140
6.9	Illustration of the $d_2$ coefficient for the nine integers . . . . .	141
6.10	A pair of attributes for the nine integers. . . . .	141
7.1	Single-linkage GDC . . . . .	149
7.2	Crime data for the Commonwealth of New Jersey . . . . .	151
7.3	Standardized crime data for New Jersey . . . . .	151
7.4	Single-linkage for NJ crimes data . . . . .	152
7.5	Complete-linkage for NJ crimes data . . . . .	152
7.6	Labels for vector representations of a DC . . . . .	152
7.7	Average Vietnam combat deaths . . . . .	165
7.8	Combat deaths in Vietnam . . . . .	167