

DATA CLUSTERING

Algorithms and
Applications

Edited by

Charu C. Aggarwal
Chandan K. Reddy



CRC Press
Taylor & Francis Group
Boca Raton London New York

CRC Press is an imprint of the
Taylor & Francis Group, an **informa** business
A CHAPMAN & HALL BOOK

CRC Press
Taylor & Francis Group
6000 Broken Sound Parkway NW, Suite 300
Boca Raton, FL 33487-2742

© 2014 by Taylor & Francis Group, LLC
CRC Press is an imprint of Taylor & Francis Group, an Informa business

No claim to original U.S. Government works

Printed on acid-free paper
Version Date: 20130508

International Standard Book Number-13: 978-1-4665-5821-2 (Hardback)

This book contains information obtained from authentic and highly regarded sources. Reasonable efforts have been made to publish reliable data and information, but the author and publisher cannot assume responsibility for the validity of all materials or the consequences of their use. The authors and publishers have attempted to trace the copyright holders of all material reproduced in this publication and apologize to copyright holders if permission to publish in this form has not been obtained. If any copyright material has not been acknowledged please write and let us know so we may rectify in any future reprint.

Except as permitted under U.S. Copyright Law, no part of this book may be reprinted, reproduced, transmitted, or utilized in any form by any electronic, mechanical, or other means, now known or hereafter invented, including photocopying, microfilming, and recording, or in any information storage or retrieval system, without written permission from the publishers.

For permission to photocopy or use material electronically from this work, please access www.copyright.com (<http://www.copyright.com/>) or contact the Copyright Clearance Center, Inc. (CCC), 222 Rosewood Drive, Danvers, MA 01923, 978-750-8400. CCC is a not-for-profit organization that provides licenses and registration for a variety of users. For organizations that have been granted a photocopy license by the CCC, a separate system of payment has been arranged.

Trademark Notice: Product or corporate names may be trademarks or registered trademarks, and are used only for identification and explanation without intent to infringe.

Library of Congress Cataloging-in-Publication Data

Data clustering : algorithms and applications / [edited by] Charu C. Aggarwal, Chandan K. Reddy.
pages cm. -- (Chapman & Hall/CRC data mining and knowledge discovery series)

Includes bibliographical references and index.

ISBN 978-1-4665-5821-2 (hardback)

1. Document clustering. 2. Cluster analysis. 3. Data mining. 4. Machine theory. 5. File organization (Computer science) I. Aggarwal, Charu C., editor of compilation. II. Reddy, Chandan K., 1980- editor of compilation.

QA278.D294 2014

519.5'35--dc23

2013008698

Visit the Taylor & Francis Web site at
<http://www.taylorandfrancis.com>

and the CRC Press Web site at
<http://www.crcpress.com>

Contents

Preface	xxi
Editor Biographies	xxiii
Contributors	xxv
1 An Introduction to Cluster Analysis	1
<i>Charu C. Aggarwal</i>	
1.1 Introduction	2
1.2 Common Techniques Used in Cluster Analysis	3
1.2.1 Feature Selection Methods	4
1.2.2 Probabilistic and Generative Models	4
1.2.3 Distance-Based Algorithms	5
1.2.4 Density- and Grid-Based Methods	7
1.2.5 Leveraging Dimensionality Reduction Methods	8
1.2.5.1 Generative Models for Dimensionality Reduction	8
1.2.5.2 Matrix Factorization and Co-Clustering	8
1.2.5.3 Spectral Methods	10
1.2.6 The High Dimensional Scenario	11
1.2.7 Scalable Techniques for Cluster Analysis	13
1.2.7.1 I/O Issues in Database Management	13
1.2.7.2 Streaming Algorithms	14
1.2.7.3 The Big Data Framework	14
1.3 Data Types Studied in Cluster Analysis	15
1.3.1 Clustering Categorical Data	15
1.3.2 Clustering Text Data	16
1.3.3 Clustering Multimedia Data	16
1.3.4 Clustering Time-Series Data	17
1.3.5 Clustering Discrete Sequences	17
1.3.6 Clustering Network Data	18
1.3.7 Clustering Uncertain Data	19
1.4 Insights Gained from Different Variations of Cluster Analysis	19
1.4.1 Visual Insights	20
1.4.2 Supervised Insights	20
1.4.3 Multiview and Ensemble-Based Insights	21
1.4.4 Validation-Based Insights	21
1.5 Discussion and Conclusions	22

2	Feature Selection for Clustering: A Review	29
	<i>Salem Alelyani, Jiliang Tang, and Huan Liu</i>	
2.1	Introduction	30
2.1.1	Data Clustering	32
2.1.2	Feature Selection	32
2.1.3	Feature Selection for Clustering	33
2.1.3.1	Filter Model	34
2.1.3.2	Wrapper Model	35
2.1.3.3	Hybrid Model	35
2.2	Feature Selection for Clustering	35
2.2.1	Algorithms for Generic Data	36
2.2.1.1	Spectral Feature Selection (SPEC)	36
2.2.1.2	Laplacian Score (LS)	36
2.2.1.3	Feature Selection for Sparse Clustering	37
2.2.1.4	Localized Feature Selection Based on Scatter Separability (LFSBSS)	38
2.2.1.5	Multicluster Feature Selection (MCFS)	39
2.2.1.6	Feature Weighting k -Means	40
2.2.2	Algorithms for Text Data	41
2.2.2.1	Term Frequency (TF)	41
2.2.2.2	Inverse Document Frequency (IDF)	42
2.2.2.3	Term Frequency-Inverse Document Frequency (TF-IDF)	42
2.2.2.4	Chi Square Statistic	42
2.2.2.5	Frequent Term-Based Text Clustering	44
2.2.2.6	Frequent Term Sequence	45
2.2.3	Algorithms for Streaming Data	47
2.2.3.1	Text Stream Clustering Based on Adaptive Feature Selection (TSC-AFS)	47
2.2.3.2	High-Dimensional Projected Stream Clustering (HPStream)	48
2.2.4	Algorithms for Linked Data	50
2.2.4.1	Challenges and Opportunities	50
2.2.4.2	LUFS: An Unsupervised Feature Selection Framework for Linked Data	51
2.2.4.3	Conclusion and Future Work for Linked Data	52
2.3	Discussions and Challenges	53
2.3.1	The Chicken or the Egg Dilemma	53
2.3.2	Model Selection: K and l	54
2.3.3	Scalability	54
2.3.4	Stability	55
3	Probabilistic Models for Clustering	61
	<i>Hongbo Deng and Jiawei Han</i>	
3.1	Introduction	61
3.2	Mixture Models	62
3.2.1	Overview	62
3.2.2	Gaussian Mixture Model	64
3.2.3	Bernoulli Mixture Model	67
3.2.4	Model Selection Criteria	68
3.3	EM Algorithm and Its Variations	69
3.3.1	The General EM Algorithm	69
3.3.2	Mixture Models Revisited	73

3.3.3	Limitations of the EM Algorithm	75
3.3.4	Applications of the EM Algorithm	76
3.4	Probabilistic Topic Models	76
3.4.1	Probabilistic Latent Semantic Analysis	77
3.4.2	Latent Dirichlet Allocation	79
3.4.3	Variations and Extensions	81
3.5	Conclusions and Summary	81
4	A Survey of Partitional and Hierarchical Clustering Algorithms	87
	<i>Chandan K. Reddy and Bhanukiran Vinzamuri</i>	
4.1	Introduction	88
4.2	Partitional Clustering Algorithms	89
4.2.1	<i>K</i> -Means Clustering	89
4.2.2	Minimization of Sum of Squared Errors	90
4.2.3	Factors Affecting <i>K</i> -Means	91
	4.2.3.1 Popular Initialization Methods	91
	4.2.3.2 Estimating the Number of Clusters	92
4.2.4	Variations of <i>K</i> -Means	93
	4.2.4.1 <i>K</i> -Medoids Clustering	93
	4.2.4.2 <i>K</i> -Medians Clustering	94
	4.2.4.3 <i>K</i> -Modes Clustering	94
	4.2.4.4 Fuzzy <i>K</i> -Means Clustering	95
	4.2.4.5 <i>X</i> -Means Clustering	95
	4.2.4.6 Intelligent <i>K</i> -Means Clustering	96
	4.2.4.7 Bisecting <i>K</i> -Means Clustering	97
	4.2.4.8 Kernel <i>K</i> -Means Clustering	97
	4.2.4.9 Mean Shift Clustering	98
	4.2.4.10 Weighted <i>K</i> -Means Clustering	98
	4.2.4.11 Genetic <i>K</i> -Means Clustering	99
4.2.5	Making <i>K</i> -Means Faster	100
4.3	Hierarchical Clustering Algorithms	100
4.3.1	Agglomerative Clustering	101
	4.3.1.1 Single and Complete Link	101
	4.3.1.2 Group Averaged and Centroid Agglomerative Clustering	102
	4.3.1.3 Ward's Criterion	103
	4.3.1.4 Agglomerative Hierarchical Clustering Algorithm	103
	4.3.1.5 Lance–Williams Dissimilarity Update Formula	103
4.3.2	Divisive Clustering	104
	4.3.2.1 Issues in Divisive Clustering	104
	4.3.2.2 Divisive Hierarchical Clustering Algorithm	105
	4.3.2.3 Minimum Spanning Tree-Based Clustering	105
4.3.3	Other Hierarchical Clustering Algorithms	106
4.4	Discussion and Summary	106
5	Density-Based Clustering	111
	<i>Martin Ester</i>	
5.1	Introduction	111
5.2	DBSCAN	113
5.3	DENCLUE	115
5.4	OPTICS	116
5.5	Other Algorithms	116

5.6	Subspace Clustering	118
5.7	Clustering Networks	120
5.8	Other Directions	123
5.9	Conclusion	124
6	Grid-Based Clustering	127
	<i>Wei Cheng, Wei Wang, and Sandra Batista</i>	
6.1	Introduction	128
6.2	The Classical Algorithms	131
6.2.1	Earliest Approaches: GRIDCLUS and BANG	131
6.2.2	STING and STING+: The Statistical Information Grid Approach	132
6.2.3	WaveCluster: Wavelets in Grid-Based Clustering	134
6.3	Adaptive Grid-Based Algorithms	135
6.3.1	AMR: Adaptive Mesh Refinement Clustering	135
6.4	Axis-Shifting Grid-Based Algorithms	136
6.4.1	NSGC: New Shifting Grid Clustering Algorithm	136
6.4.2	ADCC: Adaptable Deflect and Conquer Clustering	137
6.4.3	ASGC: Axis-Shifted Grid-Clustering	137
6.4.4	GDILC: Grid-Based Density-IsoLine Clustering Algorithm	138
6.5	High-Dimensional Algorithms	139
6.5.1	CLIQUE: The Classical High-Dimensional Algorithm	139
6.5.2	Variants of CLIQUE	140
6.5.2.1	ENCLUS: Entropy-Based Approach	140
6.5.2.2	MAFIA: Adaptive Grids in High Dimensions	141
6.5.3	OptiGrid: Density-Based Optimal Grid Partitioning	141
6.5.4	Variants of the OptiGrid Approach	143
6.5.4.1	O-Cluster: A Scalable Approach	143
6.5.4.2	CBF: Cell-Based Filtering	144
6.6	Conclusions and Summary	145
7	Nonnegative Matrix Factorizations for Clustering: A Survey	149
	<i>Tao Li and Chris Ding</i>	
7.1	Introduction	150
7.1.1	Background	150
7.1.2	NMF Formulations	151
7.2	NMF for Clustering: Theoretical Foundations	151
7.2.1	NMF and K -Means Clustering	151
7.2.2	NMF and Probabilistic Latent Semantic Indexing	152
7.2.3	NMF and Kernel K -Means and Spectral Clustering	152
7.2.4	NMF Boundedness Theorem	153
7.3	NMF Clustering Capabilities	153
7.3.1	Examples	153
7.3.2	Analysis	153
7.4	NMF Algorithms	155
7.4.1	Introduction	155
7.4.2	Algorithm Development	155
7.4.3	Practical Issues in NMF Algorithms	156
7.4.3.1	Initialization	156
7.4.3.2	Stopping Criteria	156
7.4.3.3	Objective Function vs. Clustering Performance	157
7.4.3.4	Scalability	157

7.5	NMF Related Factorizations	158
7.6	NMF for Clustering: Extensions	161
7.6.1	Co-clustering	161
7.6.2	Semisupervised Clustering	162
7.6.3	Semisupervised Co-Clustering	162
7.6.4	Consensus Clustering	163
7.6.5	Graph Clustering	164
7.6.6	Other Clustering Extensions	164
7.7	Conclusions	165
8	Spectral Clustering	177
	<i>Jialu Liu and Jiawei Han</i>	
8.1	Introduction	177
8.2	Similarity Graph	179
8.3	Unnormalized Spectral Clustering	180
8.3.1	Notation	180
8.3.2	Unnormalized Graph Laplacian	180
8.3.3	Spectrum Analysis	181
8.3.4	Unnormalized Spectral Clustering Algorithm	182
8.4	Normalized Spectral Clustering	182
8.4.1	Normalized Graph Laplacian	183
8.4.2	Spectrum Analysis	184
8.4.3	Normalized Spectral Clustering Algorithm	184
8.5	Graph Cut View	185
8.5.1	Ratio Cut Relaxation	186
8.5.2	Normalized Cut Relaxation	187
8.6	Random Walks View	188
8.7	Connection to Laplacian Eigenmap	189
8.8	Connection to Kernel k -Means and Nonnegative Matrix Factorization	191
8.9	Large Scale Spectral Clustering	192
8.10	Further Reading	194
9	Clustering High-Dimensional Data	201
	<i>Arthur Zimek</i>	
9.1	Introduction	201
9.2	The “Curse of Dimensionality”	202
9.2.1	Different Aspects of the “Curse”	202
9.2.2	Consequences	206
9.3	Clustering Tasks in Subspaces of High-Dimensional Data	206
9.3.1	Categories of Subspaces	206
9.3.1.1	Axis-Parallel Subspaces	206
9.3.1.2	Arbitrarily Oriented Subspaces	207
9.3.1.3	Special Cases	207
9.3.2	Search Spaces for the Clustering Problem	207
9.4	Fundamental Algorithmic Ideas	208
9.4.1	Clustering in Axis-Parallel Subspaces	208
9.4.1.1	Cluster Model	208
9.4.1.2	Basic Techniques	208
9.4.1.3	Clustering Algorithms	210
9.4.2	Clustering in Arbitrarily Oriented Subspaces	215
9.4.2.1	Cluster Model	215

9.4.2.2	Basic Techniques and Example Algorithms	216
9.5	Open Questions and Current Research Directions	218
9.6	Conclusion	219
10	A Survey of Stream Clustering Algorithms	231
	<i>Charu C. Aggarwal</i>	
10.1	Introduction	231
10.2	Methods Based on Partitioning Representatives	233
10.2.1	The STREAM Algorithm	233
10.2.2	CluStream: The Microclustering Framework	235
10.2.2.1	Microcluster Definition	235
10.2.2.2	Pyramidal Time Frame	236
10.2.2.3	Online Clustering with CluStream	237
10.3	Density-Based Stream Clustering	239
10.3.1	DenStream: Density-Based Microclustering	240
10.3.2	Grid-Based Streaming Algorithms	241
10.3.2.1	D-Stream Algorithm	241
10.3.2.2	Other Grid-Based Algorithms	242
10.4	Probabilistic Streaming Algorithms	243
10.5	Clustering High-Dimensional Streams	243
10.5.1	The HPSTREAM Method	244
10.5.2	Other High-Dimensional Streaming Algorithms	244
10.6	Clustering Discrete and Categorical Streams	245
10.6.1	Clustering Binary Data Streams with k -Means	245
10.6.2	The StreamCluCD Algorithm	245
10.6.3	Massive-Domain Clustering	246
10.7	Text Stream Clustering	249
10.8	Other Scenarios for Stream Clustering	252
10.8.1	Clustering Uncertain Data Streams	253
10.8.2	Clustering Graph Streams	253
10.8.3	Distributed Clustering of Data Streams	254
10.9	Discussion and Conclusions	254
11	Big Data Clustering	259
	<i>Hanghang Tong and U Kang</i>	
11.1	Introduction	259
11.2	One-Pass Clustering Algorithms	260
11.2.1	CLARANS: Fighting with Exponential Search Space	260
11.2.2	BIRCH: Fighting with Limited Memory	261
11.2.3	CURE: Fighting with the Irregular Clusters	263
11.3	Randomized Techniques for Clustering Algorithms	263
11.3.1	Locality-Preserving Projection	264
11.3.2	Global Projection	266
11.4	Parallel and Distributed Clustering Algorithms	268
11.4.1	General Framework	268
11.4.2	DBDC: Density-Based Clustering	269
11.4.3	ParMETIS: Graph Partitioning	269
11.4.4	PKMeans: K -Means with MapReduce	270
11.4.5	DisCo: Co-Clustering with MapReduce	271
11.4.6	BoW: Subspace Clustering with MapReduce	272
11.5	Conclusion	274

12 Clustering Categorical Data	277
<i>Bill Andreopoulos</i>	
12.1 Introduction	278
12.2 Goals of Categorical Clustering	279
12.2.1 Clustering Road Map	280
12.3 Similarity Measures for Categorical Data	282
12.3.1 The Hamming Distance in Categorical and Binary Data	282
12.3.2 Probabilistic Measures	283
12.3.3 Information-Theoretic Measures	283
12.3.4 Context-Based Similarity Measures	284
12.4 Descriptions of Algorithms	284
12.4.1 Partition-Based Clustering	284
12.4.1.1 k -Modes	284
12.4.1.2 k -Prototypes (Mixed Categorical and Numerical)	285
12.4.1.3 Fuzzy k -Modes	286
12.4.1.4 Squeezer	286
12.4.1.5 COOLCAT	286
12.4.2 Hierarchical Clustering	287
12.4.2.1 ROCK	287
12.4.2.2 COBWEB	288
12.4.2.3 LIMBO	289
12.4.3 Density-Based Clustering	289
12.4.3.1 Projected (Subspace) Clustering	290
12.4.3.2 CACTUS	290
12.4.3.3 CLICKS	291
12.4.3.4 STIRR	291
12.4.3.5 CLOPE	292
12.4.3.6 HIERDENC: Hierarchical Density-Based Clustering	292
12.4.3.7 MULIC: Multiple Layer Incremental Clustering	293
12.4.4 Model-Based Clustering	296
12.4.4.1 BILCOM Empirical Bayesian (Mixed Categorical and Numerical)	296
12.4.4.2 AutoClass (Mixed Categorical and Numerical)	296
12.4.4.3 SVM Clustering (Mixed Categorical and Numerical)	297
12.5 Conclusion	298
13 Document Clustering: The Next Frontier	305
<i>David C. Anastasiu, Andrea Tagarelli, and George Karypis</i>	
13.1 Introduction	306
13.2 Modeling a Document	306
13.2.1 Preliminaries	306
13.2.2 The Vector Space Model	307
13.2.3 Alternate Document Models	309
13.2.4 Dimensionality Reduction for Text	309
13.2.5 Characterizing Extremes	310
13.3 General Purpose Document Clustering	311
13.3.1 Similarity/Dissimilarity-Based Algorithms	311
13.3.2 Density-Based Algorithms	312
13.3.3 Adjacency-Based Algorithms	313
13.3.4 Generative Algorithms	313
13.4 Clustering Long Documents	315

13.4.1	Document Segmentation	315
13.4.2	Clustering Segmented Documents	317
13.4.3	Simultaneous Segment Identification and Clustering	321
13.5	Clustering Short Documents	323
13.5.1	General Methods for Short Document Clustering	323
13.5.2	Clustering with Knowledge Infusion	324
13.5.3	Clustering Web Snippets	325
13.5.4	Clustering Microblogs	326
13.6	Conclusion	328
14	Clustering Multimedia Data	339
	<i>Shen-Fu Tsai, Guo-Jun Qi, Shiyu Chang, Min-Hsuan Tsai, and Thomas S. Huang</i>	
14.1	Introduction	340
14.2	Clustering with Image Data	340
14.2.1	Visual Words Learning	341
14.2.2	Face Clustering and Annotation	342
14.2.3	Photo Album Event Recognition	343
14.2.4	Image Segmentation	344
14.2.5	Large-Scale Image Classification	345
14.3	Clustering with Video and Audio Data	347
14.3.1	Video Summarization	348
14.3.2	Video Event Detection	349
14.3.3	Video Story Clustering	350
14.3.4	Music Summarization	350
14.4	Clustering with Multimodal Data	351
14.5	Summary and Future Directions	353
15	Time-Series Data Clustering	357
	<i>Dimitrios Kotsakos, Goce Trajcevski, Dimitrios Gunopulos, and Charu C. Aggarwal</i>	
15.1	Introduction	358
15.2	The Diverse Formulations for Time-Series Clustering	359
15.3	Online Correlation-Based Clustering	360
15.3.1	Selective Muscles and Related Methods	361
15.3.2	Sensor Selection Algorithms for Correlation Clustering	362
15.4	Similarity and Distance Measures	363
15.4.1	Univariate Distance Measures	363
15.4.1.1	L_p Distance	363
15.4.1.2	Dynamic Time Warping Distance	364
15.4.1.3	EDIT Distance	365
15.4.1.4	Longest Common Subsequence	365
15.4.2	Multivariate Distance Measures	366
15.4.2.1	Multidimensional L_p Distance	366
15.4.2.2	Multidimensional DTW	367
15.4.2.3	Multidimensional LCSS	368
15.4.2.4	Multidimensional Edit Distance	368
15.4.2.5	Multidimensional Subsequence Matching	368
15.5	Shape-Based Time-Series Clustering Techniques	369
15.5.1	k -Means Clustering	370
15.5.2	Hierarchical Clustering	371
15.5.3	Density-Based Clustering	372

15.5.4	Trajectory Clustering	372
15.6	Time-Series Clustering Applications	374
15.7	Conclusions	375
16	Clustering Biological Data	381
	<i>Chandan K. Reddy, Mohammad Al Hasan, and Mohammed J. Zaki</i>	
16.1	Introduction	382
16.2	Clustering Microarray Data	383
16.2.1	Proximity Measures	383
16.2.2	Categorization of Algorithms	384
16.2.3	Standard Clustering Algorithms	385
16.2.3.1	Hierarchical Clustering	385
16.2.3.2	Probabilistic Clustering	386
16.2.3.3	Graph-Theoretic Clustering	386
16.2.3.4	Self-Organizing Maps	387
16.2.3.5	Other Clustering Methods	387
16.2.4	Biclustering	388
16.2.4.1	Types and Structures of Biclusters	389
16.2.4.2	Biclustering Algorithms	390
16.2.4.3	Recent Developments	391
16.2.5	Triclustering	391
16.2.6	Time-Series Gene Expression Data Clustering	392
16.2.7	Cluster Validation	393
16.3	Clustering Biological Networks	394
16.3.1	Characteristics of PPI Network Data	394
16.3.2	Network Clustering Algorithms	394
16.3.2.1	Molecular Complex Detection	394
16.3.2.2	Markov Clustering	395
16.3.2.3	Neighborhood Search Methods	395
16.3.2.4	Clique Percolation Method	395
16.3.2.5	Ensemble Clustering	396
16.3.2.6	Other Clustering Methods	396
16.3.3	Cluster Validation and Challenges	397
16.4	Biological Sequence Clustering	397
16.4.1	Sequence Similarity Metrics	397
16.4.1.1	Alignment-Based Similarity	398
16.4.1.2	Keyword-Based Similarity	398
16.4.1.3	Kernel-Based Similarity	399
16.4.1.4	Model-Based Similarity	399
16.4.2	Sequence Clustering Algorithms	399
16.4.2.1	Subsequence-Based Clustering	399
16.4.2.2	Graph-Based Clustering	400
16.4.2.3	Probabilistic Models	402
16.4.2.4	Suffix Tree and Suffix Array-Based Method	403
16.5	Software Packages	403
16.6	Discussion and Summary	405

17 Network Clustering	415
<i>Srinivasan Parthasarathy and S M Faisal</i>	
17.1 Introduction	416
17.2 Background and Nomenclature	417
17.3 Problem Definition	417
17.4 Common Evaluation Criteria	418
17.5 Partitioning with Geometric Information	419
17.5.1 Coordinate Bisection	419
17.5.2 Inertial Bisection	419
17.5.3 Geometric Partitioning	420
17.6 Graph Growing and Greedy Algorithms	421
17.6.1 Kernighan-Lin Algorithm	422
17.7 Agglomerative and Divisive Clustering	423
17.8 Spectral Clustering	424
17.8.1 Similarity Graphs	425
17.8.2 Types of Similarity Graphs	425
17.8.3 Graph Laplacians	426
17.8.3.1 Unnormalized Graph Laplacian	426
17.8.3.2 Normalized Graph Laplacians	427
17.8.4 Spectral Clustering Algorithms	427
17.9 Markov Clustering	428
17.9.1 Regularized MCL (RMCL): Improvement over MCL	429
17.10 Multilevel Partitioning	430
17.11 Local Partitioning Algorithms	432
17.12 Hypergraph Partitioning	433
17.13 Emerging Methods for Partitioning Special Graphs	435
17.13.1 Bipartite Graphs	435
17.13.2 Dynamic Graphs	436
17.13.3 Heterogeneous Networks	437
17.13.4 Directed Networks	438
17.13.5 Combining Content and Relationship Information	439
17.13.6 Networks with Overlapping Communities	440
17.13.7 Probabilistic Methods	442
17.14 Conclusion	443
18 A Survey of Uncertain Data Clustering Algorithms	457
<i>Charu C. Aggarwal</i>	
18.1 Introduction	457
18.2 Mixture Model Clustering of Uncertain Data	459
18.3 Density-Based Clustering Algorithms	460
18.3.1 FDBSCAN Algorithm	460
18.3.2 FOPTICS Algorithm	461
18.4 Partitional Clustering Algorithms	462
18.4.1 The UK-Means Algorithm	462
18.4.2 The CK-Means Algorithm	463
18.4.3 Clustering Uncertain Data with Voronoi Diagrams	464
18.4.4 Approximation Algorithms for Clustering Uncertain Data	464
18.4.5 Speeding Up Distance Computations	465
18.5 Clustering Uncertain Data Streams	466
18.5.1 The UMicro Algorithm	466
18.5.2 The LuMicro Algorithm	471

18.5.3	Enhancements to Stream Clustering	471
18.6	Clustering Uncertain Data in High Dimensionality	472
18.6.1	Subspace Clustering of Uncertain Data	473
18.6.2	UPStream: Projected Clustering of Uncertain Data Streams	474
18.7	Clustering with the Possible Worlds Model	477
18.8	Clustering Uncertain Graphs	478
18.9	Conclusions and Summary	478
19	Concepts of Visual and Interactive Clustering	483
	<i>Alexander Hinneburg</i>	
19.1	Introduction	483
19.2	Direct Visual and Interactive Clustering	484
19.2.1	Scatterplots	485
19.2.2	Parallel Coordinates	488
19.2.3	Discussion	491
19.3	Visual Interactive Steering of Clustering	491
19.3.1	Visual Assessment of Convergence of Clustering Algorithm	491
19.3.2	Interactive Hierarchical Clustering	492
19.3.3	Visual Clustering with SOMs	494
19.3.4	Discussion	494
19.4	Interactive Comparison and Combination of Clusterings	495
19.4.1	Space of Clusterings	495
19.4.2	Visualization	497
19.4.3	Discussion	497
19.5	Visualization of Clusters for Sense-Making	497
19.6	Summary	500
20	Semisupervised Clustering	505
	<i>Amrudin Agovic and Arindam Banerjee</i>	
20.1	Introduction	506
20.2	Clustering with Pointwise and Pairwise Semisupervision	507
20.2.1	Semisupervised Clustering Based on Seeding	507
20.2.2	Semisupervised Clustering Based on Pairwise Constraints	508
20.2.3	Active Learning for Semisupervised Clustering	511
20.2.4	Semisupervised Clustering Based on User Feedback	512
20.2.5	Semisupervised Clustering Based on Nonnegative Matrix Factorization	513
20.3	Semisupervised Graph Cuts	513
20.3.1	Semisupervised Unnormalized Cut	515
20.3.2	Semisupervised Ratio Cut	515
20.3.3	Semisupervised Normalized Cut	516
20.4	A Unified View of Label Propagation	517
20.4.1	Generalized Label Propagation	517
20.4.2	Gaussian Fields	517
20.4.3	Tikhonov Regularization (TIKREG)	518
20.4.4	Local and Global Consistency	518
20.4.5	Related Methods	519
20.4.5.1	Cluster Kernels	519
20.4.5.2	Gaussian Random Walks EM (GWEM)	519
20.4.5.3	Linear Neighborhood Propagation	520
20.4.6	Label Propagation and Green's Function	521
20.4.7	Label Propagation and Semisupervised Graph Cuts	521

20.5	Semisupervised Embedding	521
20.5.1	Nonlinear Manifold Embedding	522
20.5.2	Semisupervised Embedding	522
20.5.2.1	Unconstrained Semisupervised Embedding	523
20.5.2.2	Constrained Semisupervised Embedding	523
20.6	Comparative Experimental Analysis	524
20.6.1	Experimental Results	524
20.6.2	Semisupervised Embedding Methods	529
20.7	Conclusions	530
21	Alternative Clustering Analysis: A Review	535
	<i>James Bailey</i>	
21.1	Introduction	535
21.2	Technical Preliminaries	537
21.3	Multiple Clustering Analysis Using Alternative Clusterings	538
21.3.1	Alternative Clustering Algorithms: A Taxonomy	538
21.3.2	Unguided Generation	539
21.3.2.1	Naive	539
21.3.2.2	Meta Clustering	539
21.3.2.3	Eigenvectors of the Laplacian Matrix	540
21.3.2.4	Decorrelated k -Means and Convolutional EM	540
21.3.2.5	CAMI	540
21.3.3	Guided Generation with Constraints	541
21.3.3.1	COALA	541
21.3.3.2	Constrained Optimization Approach	541
21.3.3.3	MAXIMUS	542
21.3.4	Orthogonal Transformation Approaches	543
21.3.4.1	Orthogonal Views	543
21.3.4.2	ADFT	543
21.3.5	Information Theoretic	544
21.3.5.1	Conditional Information Bottleneck (CIB)	544
21.3.5.2	Conditional Ensemble Clustering	544
21.3.5.3	NACI	544
21.3.5.4	mSC	545
21.4	Connections to Multiview Clustering and Subspace Clustering	545
21.5	Future Research Issues	547
21.6	Summary	547
22	Cluster Ensembles: Theory and Applications	551
	<i>Joydeep Ghosh and Ayan Acharya</i>	
22.1	Introduction	551
22.2	The Cluster Ensemble Problem	554
22.3	Measuring Similarity Between Clustering Solutions	555
22.4	Cluster Ensemble Algorithms	558
22.4.1	Probabilistic Approaches to Cluster Ensembles	558
22.4.1.1	A Mixture Model for Cluster Ensembles (MMCE)	558
22.4.1.2	Bayesian Cluster Ensembles (BCE)	558
22.4.1.3	Nonparametric Bayesian Cluster Ensembles (NPBCE)	559
22.4.2	Pairwise Similarity-Based Approaches	560
22.4.2.1	Methods Based on Ensemble Co-Association Matrix	560

22.4.2.2	Relating Consensus Clustering to Other Optimization Formulations	562
22.4.3	Direct Approaches Using Cluster Labels	562
22.4.3.1	Graph Partitioning	562
22.4.3.2	Cumulative Voting	563
22.5	Applications of Consensus Clustering	564
22.5.1	Gene Expression Data Analysis	564
22.5.2	Image Segmentation	564
22.6	Concluding Remarks	566
23	Clustering Validation Measures	571
	<i>Hui Xiong and Zhongmou Li</i>	
23.1	Introduction	572
23.2	External Clustering Validation Measures	573
23.2.1	An Overview of External Clustering Validation Measures	574
23.2.2	Defective Validation Measures	575
23.2.2.1	<i>K</i> -Means: The Uniform Effect	575
23.2.2.2	A Necessary Selection Criterion	576
23.2.2.3	The Cluster Validation Results	576
23.2.2.4	The Issues with the Defective Measures	577
23.2.2.5	Improving the Defective Measures	577
23.2.3	Measure Normalization	577
23.2.3.1	Normalizing the Measures	578
23.2.3.2	The <i>DCV</i> Criterion	581
23.2.3.3	The Effect of Normalization	583
23.2.4	Measure Properties	584
23.2.4.1	The Consistency Between Measures	584
23.2.4.2	Properties of Measures	586
23.2.4.3	Discussions	589
23.3	Internal Clustering Validation Measures	589
23.3.1	An Overview of Internal Clustering Validation Measures	589
23.3.2	Understanding of Internal Clustering Validation Measures	592
23.3.2.1	The Impact of Monotonicity	592
23.3.2.2	The Impact of Noise	593
23.3.2.3	The Impact of Density	594
23.3.2.4	The Impact of Subclusters	595
23.3.2.5	The Impact of Skewed Distributions	596
23.3.2.6	The Impact of Arbitrary Shapes	598
23.3.3	Properties of Measures	600
23.4	Summary	601
24	Educational and Software Resources for Data Clustering	607
	<i>Charu C. Aggarwal and Chandan K. Reddy</i>	
24.1	Introduction	607
24.2	Educational Resources	608
24.2.1	Books on Data Clustering	608
24.2.2	Popular Survey Papers on Data Clustering	608
24.3	Software for Data Clustering	610
24.3.1	Free and Open-Source Software	610
24.3.1.1	General Clustering Software	610
24.3.1.2	Specialized Clustering Software	610

24.3.2	Commercial Packages	611
24.3.3	Data Benchmarks for Software and Research	611
24.4	Summary	612
Index		617

Preface

The problem of clustering is perhaps one of the most widely studied in the data mining and machine learning communities. This problem has been studied by researchers from several disciplines over five decades. Applications of clustering include a wide variety of problem domains such as text, multimedia, social networks, and biological data. Furthermore, the problem may be encountered in a number of different scenarios such as streaming or uncertain data. Clustering is a rather diverse topic, and the underlying algorithms depend greatly on the data domain and problem scenario.

Therefore, this book will focus on three primary aspects of data clustering. The first set of chapters will focus on the core methods for data clustering. These include methods such as probabilistic clustering, density-based clustering, grid-based clustering, and spectral clustering. The second set of chapters will focus on different problem domains and scenarios such as multimedia data, text data, biological data, categorical data, network data, data streams and uncertain data. The third set of chapters will focus on different detailed insights from the clustering process, because of the subjectivity of the clustering process, and the many different ways in which the same data set can be clustered. How do we know that a particular clustering is good or that it solves the needs of the application? There are numerous ways in which these issues can be explored. The exploration could be through interactive visualization and human interaction, external knowledge-based supervision, explicitly examining the multiple solutions in order to evaluate different possibilities, combining the multiple solutions in order to create more robust ensembles, or trying to judge the quality of different solutions with the use of different validation criteria.

The clustering problem has been addressed by a number of different communities such as pattern recognition, databases, data mining and machine learning. In some cases, the work by the different communities tends to be fragmented and has not been addressed in a unified way. This book will make a conscious effort to address the work of the different communities in a unified way. The book will start off with an overview of the basic methods in data clustering, and then discuss progressively more refined and complex methods for data clustering. Special attention will also be paid to more recent problem domains such as graphs and social networks.

The chapters in the book will be divided into three types:

- **Method Chapters:** These chapters discuss the *key techniques* which are commonly used for clustering such as feature selection, agglomerative clustering, partitional clustering, density-based clustering, probabilistic clustering, grid-based clustering, spectral clustering, and non-negative matrix factorization.
- **Domain Chapters:** These chapters discuss the specific methods used for different *domains* of data such as categorical data, text data, multimedia data, graph data, biological data, stream data, uncertain data, time series clustering, high-dimensional clustering, and big data. Many of these chapters can also be considered application chapters, because they explore the specific characteristics of the problem in a particular domain.
- **Variations and Insights:** These chapters discuss the *key variations* on the clustering process such as semi-supervised clustering, interactive clustering, multi-view clustering, cluster ensembles, and cluster validation. Such methods are typically used in order to obtain detailed insights from the clustering process, and also to explore different possibilities on the clustering process through either supervision, human intervention, or through automated generation

of alternative clusters. The methods for cluster validation also provide an idea of the quality of the underlying clusters.

This book is designed to be comprehensive in its coverage of the entire area of clustering, and it is hoped that it will serve as a knowledgeable compendium to students and researchers.

Editor Biographies

Charu C. Aggarwal is a Research Scientist at the IBM T. J. Watson Research Center in Yorktown Heights, New York. He completed his B.S. from IIT Kanpur in 1993 and his Ph.D. from Massachusetts Institute of Technology in 1996. His research interest during his Ph.D. years was in combinatorial optimization (network flow algorithms), and his thesis advisor was Professor James B. Orlin. He has since worked in the field of performance analysis, databases, and data mining. He has published over 200 papers in refereed conferences and journals, and has applied for or been granted over 80 patents. He is author or editor of nine books, including this one. Because of the commercial value of the above-mentioned patents, he has received several invention achievement awards and has thrice been designated a Master Inventor at IBM. He is a recipient of an IBM Corporate Award (2003) for his work on bioterrorist threat detection in data streams, a recipient of the IBM Outstanding Innovation Award (2008) for his scientific contributions to privacy technology, and a recipient of an IBM Research Division Award (2008) for his scientific contributions to data stream research. He has served on the program committees of most major database/data mining conferences, and served as program vice-chairs of the SIAM Conference on Data Mining (2007), the IEEE ICDM Conference (2007), the WWW Conference (2009), and the IEEE ICDM Conference (2009). He served as an associate editor of the *IEEE Transactions on Knowledge and Data Engineering Journal* from 2004 to 2008. He is an associate editor of the *ACM TKDD Journal*, an action editor of the *Data Mining and Knowledge Discovery Journal*, an associate editor of *ACM SIGKDD Explorations*, and an associate editor of the *Knowledge and Information Systems Journal*. He is a fellow of the IEEE for “contributions to knowledge discovery and data mining techniques”, and a life-member of the ACM.

Chandan K. Reddy is an Assistant Professor in the Department of Computer Science at Wayne State University. He received his Ph.D. from Cornell University and M.S. from Michigan State University. His primary research interests are in the areas of data mining and machine learning with applications to healthcare, bioinformatics, and social network analysis. His research is funded by the National Science Foundation, the National Institutes of Health, Department of Transportation, and the Susan G. Komen for the Cure Foundation. He has published over 40 peer-reviewed articles in leading conferences and journals. He received the Best Application Paper Award at the ACM SIGKDD conference in 2010 and was a finalist of the INFORMS Franz Edelman Award Competition in 2011. He is a member of IEEE, ACM, and SIAM.



Contributors

Ayan Acharya

University of Texas
Austin, Texas

Charu C. Aggarwal

IBM T. J. Watson Research Center
Yorktown Heights, New York

Amrudin Agovic

Reliancy, LLC
Saint Louis Park, Minnesota

Mohammad Al Hasan

Indiana University - Purdue University
Indianapolis, Indiana

Salem Alelyani

Arizona State University
Tempe, Arizona

David C. Anastasiu

University of Minnesota
Minneapolis, Minnesota

Bill Andreopoulos

Lawrence Berkeley National Laboratory
Berkeley, California

James Bailey

The University of Melbourne
Melbourne, Australia

Arindam Banerjee

University of Minnesota
Minneapolis, Minnesota

Sandra Batista

Duke University
Durham, North Carolina

Shiyu Chang

University of Illinois at Urbana-Champaign
Urbana, Illinois

Wei Cheng

University of North Carolina
Chapel Hill, North Carolina

Hongbo Deng

University of Illinois at Urbana-Champaign
Urbana, Illinois

Cha-charis Ding

University of Texas
Arlington, Texas

Martin Ester

Simon Fraser University
British Columbia, Canada

S M Faisal

The Ohio State University
Columbus, Ohio

Joydeep Ghosh

University of Texas
Austin, Texas

Dimitrios Gunopulos

University of Athens
Athens, Greece

Jiawei Han

University of Illinois at Urbana-Champaign
Urbana, Illinois

Alexander Hinneburg

Martin-Luther University
Halle/Saale, Germany

Thomas S. Huang

University of Illinois at Urbana-Champaign
Urbana, Illinois

U Kang

KAIST
Seoul, Korea

George Karypis

University of Minnesota
Minneapolis, Minnesota

Dimitrios Kotsakos

University of Athens
Athens, Greece

Tao Li

Florida International University
Miami, Florida

Zhongmou Li

Rutgers University
New Brunswick, New Jersey

Huan Liu

Arizona State University
Tempe, Arizona

Jialu Liu

University of Illinois at Urbana-Champaign
Urbana, Illinois

Srinivasan Parthasarathy

The Ohio State University
Columbus, Ohio

Guo-Jun Qi

University of Illinois at Urbana-Champaign
Urbana, Illinois

Chandan K. Reddy

Wayne State University
Detroit, Michigan

Andrea Tagarelli

University of Calabria
Arcavacata di Rende, Italy

Jiliang Tang

Arizona State University
Tempe, Arizona

Hanghang Tong

IBM T. J. Watson Research Center
Yorktown Heights, New York

Goce Trajcevski

Northwestern University
Evanston, Illinois

Min-Hsuan Tsai

University of Illinois at Urbana-Champaign
Urbana, Illinois

Shen-Fu Tsai

Microsoft Inc.
Redmond, Washington

Bhanukiran Vinzamuri

Wayne State University
Detroit, Michigan

Wei Wang

University of California
Los Angeles, California

Hui Xiong

Rutgers University
New Brunswick, New Jersey

Mohammed J. Zaki

Rensselaer Polytechnic Institute
Troy, New York

Arthur Zimek

University of Alberta
Edmonton, Canada

Chapter 1

An Introduction to Cluster Analysis

Charu C. Aggarwal

IBM T. J. Watson Research Center

Yorktown Heights, NY

charu@us.ibm.com

1.1	Introduction	1
1.2	Common Techniques Used in Cluster Analysis	3
1.2.1	Feature Selection Methods	4
1.2.2	Probabilistic and Generative Models	4
1.2.3	Distance-Based Algorithms	5
1.2.4	Density- and Grid-Based Methods	7
1.2.5	Leveraging Dimensionality Reduction Methods	8
1.2.5.1	Generative Models for Dimensionality Reduction	8
1.2.5.2	Matrix Factorization and Co-Clustering	8
1.2.5.3	Spectral Methods	10
1.2.6	The High Dimensional Scenario	11
1.2.7	Scalable Techniques for Cluster Analysis	13
1.2.7.1	I/O Issues in Database Management	13
1.2.7.2	Streaming Algorithms	14
1.2.7.3	The Big Data Framework	14
1.3	Data Types Studied in Cluster Analysis	15
1.3.1	Clustering Categorical Data	15
1.3.2	Clustering Text Data	16
1.3.3	Clustering Multimedia Data	16
1.3.4	Clustering Time-Series Data	17
1.3.5	Clustering Discrete Sequences	17
1.3.6	Clustering Network Data	18
1.3.7	Clustering Uncertain Data	19
1.4	Insights Gained from Different Variations of Cluster Analysis	19
1.4.1	Visual Insights	20
1.4.2	Supervised Insights	20
1.4.3	Multiview and Ensemble-Based Insights	21
1.4.4	Validation-Based Insights	21
1.5	Discussion and Conclusions	22
	Bibliography	23

1.1 Introduction

The problem of data clustering has been widely studied in the data mining and machine learning literature because of its numerous applications to summarization, learning, segmentation, and target marketing [46, 47, 52]. In the absence of specific labeled information, clustering can be considered a concise model of the data which can be interpreted in the sense of either a summary or a generative model. The basic problem of clustering may be stated as follows:

Given a set of data points, partition them into a set of groups which are as similar as possible.

Note that this is a very rough definition, and the variations in the problem definition can be significant, depending upon the specific model used. For example, a generative model may define similarity on the basis of a probabilistic generative mechanism, whereas a distance-based approach will use a traditional distance function for quantification. Furthermore, the specific data type also has a significant impact on the problem definition.

Some common application domains in which the clustering problem arises are as follows:

- **Intermediate Step for other fundamental data mining problems:** Since a clustering can be considered a form of data summarization, it often serves as a key intermediate step for many fundamental data mining problems such as classification or outlier analysis. A compact summary of the data is often useful for different kinds of application-specific insights.
- **Collaborative Filtering:** In collaborative filtering methods, the clustering provides a summarization of like-minded users. The ratings provided by the different users for each other are used in order to perform the collaborative filtering. This can be used to provide recommendations in a variety of applications.
- **Customer Segmentation:** This application is quite similar to collaborative filtering, since it creates groups of similar customers in the data. The major difference from collaborative filtering is that instead of using rating information, arbitrary attributes about the objects may be used for clustering purposes.
- **Data Summarization:** Many clustering methods are closely related to dimensionality reduction methods. Such methods can be considered a form of data summarization. Data summarization can be helpful in creating compact data representations, which are easier to process and interpret in a wide variety of applications.
- **Dynamic Trend Detection:** Many forms of dynamic and streaming algorithms can be used to perform trend detection in a wide variety of social networking applications. In such applications, the data is dynamically clustered in a streaming fashion and can be used in order to determine important patterns of changes. Examples of such streaming data could be multidimensional data, text streams, streaming time-series data, and trajectory data. Key trends and events in the data can be discovered with the use of clustering methods.
- **Multimedia Data Analysis:** A variety of different kinds of documents such as images, audio or video, fall in the general category of multimedia data. The determination of similar segments has numerous applications, such as the determination of similar snippets of music or similar photographs. In many cases, the data may be multimodal and may contain different types. In such cases, the problem becomes even more challenging.
- **Biological Data Analysis:** Biological data has become pervasive in the last few years, because of the success of the human genome effort and the increasing ability to collect different kinds of gene expression data. Biological data is usually structured either as sequences or as

networks. Clustering algorithms provide good ideas of the key trends in the data, as well as the unusual sequences.

- **Social Network Analysis:** In these applications, the structure of a social network is used in order to determine the important communities in the underlying network. Community detection has important applications in social network analysis, because it provides an important understanding of the community structure in the network. Clustering also has applications to social network summarization, which is useful in a number of applications.

The aforementioned list of applications is not exhaustive by any means; nevertheless it represents a good cross-section of the wide diversity of problems which can be addressed with clustering algorithms.

The work in the data clustering area typically falls into a number of broad categories;

- **Technique-centered:** Since clustering is a rather popular problem, it is not surprising that numerous methods, such as probabilistic techniques, distance-based techniques, spectral techniques, density-based techniques, and dimensionality-reduction based techniques, are used for the clustering process. Each of these methods has its own advantages and disadvantages, and may work well in different scenarios and problem domains. Certain kinds of data types such as high dimensional data, big data, or streaming data have their own set of challenges and often require specialized techniques.
- **Data-Type Centered:** Different applications create different kinds of data types with different properties. For example, an ECG machine will produce time series data points which are highly correlated with one another, whereas a social network will generate a mixture of document and structural data. Some of the most common examples are categorical data, time series data, discrete sequences, network data, and probabilistic data. Clearly, the nature of the data greatly impacts the choice of methodology used for the clustering process. Furthermore, some data types are more difficult than others because of the separation between different kinds of attributes such as behavior or contextual attributes.
- **Additional Insights from Clustering Variations:** A number of insights have also been designed for different kinds of clustering variations. For example, visual analysis, supervised analysis, ensemble-analysis, or multiview analysis can be used in order to gain additional insights. Furthermore, the issue of cluster validation is also important from the perspective of gaining specific insights about the performance of the clustering.

This chapter will discuss each of these issues in detail, and will also discuss how the organization of the book relates to these different areas of clustering. The chapter is organized as follows. The next section discusses the common techniques which are used in cluster analysis. Section 1.3 explores the use of different data types in the clustering process. Section 1.4 discusses the use of different variations of data clustering. Section 1.5 offers the conclusions and summary.

1.2 Common Techniques Used in Cluster Analysis

The clustering problems can be addressed using a wide variation of methods. In addition, the data preprocessing phase requires dedicated techniques of its own. A number of good books and surveys discuss these issues [14, 20, 31, 37, 46, 47, 48, 52, 65, 80, 81]. The most common techniques which are used for clustering are discussed in this section.

1.2.1 Feature Selection Methods

The feature selection phase is an important preprocessing step which is needed in order to enhance the quality of the underlying clustering. Not all features are equally relevant to finding the clusters, since some may be more noisy than other. Therefore, it is often helpful to utilize a preprocessing phase in which the noisy and irrelevant features are pruned from contention. Feature selection and dimensionality reduction are closely related. In feature selection, original subsets of the features are selected. In dimensionality reduction, linear combinations of features may be used in techniques such as principal component analysis [50] in order to further enhance the feature selection effect. The advantage of the former is greater interpretability, whereas the advantage of the latter is that a lesser number of transformed directions is required for the representation process. Chapter 2 of this book will discuss such feature selection methods in detail. A comprehensive book on feature selection may be found in [61].

It should be noted that feature selection can also be integrated directly into the clustering algorithm to gain better locality specific insights. This is particularly useful, when different features are relevant to different localities of the data. The motivating factor for high dimensional subspace clustering algorithms is the failure of global feature selection algorithms. As noted in [9]: “... in many real data examples, some points are correlated with respect to a given set of dimensions and others are correlated with respect to different dimensions. Thus, it may not always be feasible to prune off too many dimensions without at the same time incurring a substantial loss of information” p.61. Therefore, the use of local feature selection, by integrating the feature selection process into the algorithm, is the best way of achieving this goal. Such local feature selections can also be extended to the dimensionality reduction problem [8, 19] and are sometimes referred to as *local dimensionality reduction*. Such methods are discussed in detail in Chapter 9. Furthermore, Chapter 2 also discusses the connections of these classes of methods to the problem of feature selection.

1.2.2 Probabilistic and Generative Models

In probabilistic models, the core idea is to model the data from a *generative process*. First, a specific form of the generative model (e.g., mixture of Gaussians) is assumed, and then the parameters of this model are estimated with the use of the Expectation Maximization (EM) algorithm [27]. The available data set is used to estimate the parameters in such a way that they have a *maximum likelihood fit* to the generative model. Given this model, we then estimate the generative probabilities (or fit probabilities) of the underlying data points. Data points which fit the distribution well will have high fit probabilities, whereas anomalies will have very low fit probabilities.

The broad principle of a mixture-based generative model is to *assume* that the data were generated from a mixture of k distributions with the probability distributions $G_1 \dots G_k$ with the use of the following process:

- Pick a data distribution with prior probability α_i , where $i \in \{1 \dots k\}$, in order to pick one of the k distributions. Let us assume that the r th one is picked.
- Generate a data point from G_r .

The probability distribution G_r is picked from a host of different possibilities. Note that this generative process requires the determination of several parameters such as the prior probabilities and the model parameters for each distribution G_r . Models with different levels of flexibility may be designed depending upon whether the prior probabilities are specified as part of the problem setting, or whether interattribute correlations are assumed within a component of the mixture. Note that the model parameters and the probability of assignment of data points to clusters are dependent on one another in a circular way. Iterative methods are therefore desirable in order to resolve this circularity. The generative models are typically solved with the use of an EM approach, which starts off with

a random or heuristic initialization and then iteratively uses two steps to resolve the circularity in computation:

- (E-Step) Determine the expected probability of assignment of data points to clusters with the use of current model parameters.
- (M-Step) Determine the optimum model parameters of each mixture by using the assignment probabilities as weights.

One nice property of EM-models is that they can be generalized relatively easily to different kinds of data, as long as the generative model for each component is properly selected for the individual mixture component G_r . Some examples are as follows:

- For numerical data, a Gaussian mixture model may be used in order to model each component G_r . Such a model is discussed in detail in Chapter 3.
- For categorical data, a Bernoulli model may be used for G_r in order to model the generation of the discrete values.
- For sequence data, a Hidden Markov Model (HMM) may be used for G_r in order to model the generation of a sequence. Interestingly, an HMM is itself a special kind of mixture model in which the different components of the mixture are dependent on each other through transitions. Thus, the clustering of sequence data with a mixture of HMMs can be considered a two-level mixture model.

Generative models are among the most fundamental of all clustering methods, because they try to understand the underlying *process* through which a cluster is generated. A number of interesting connections exist between other clustering methods and generative models, by considering special cases in terms of prior probabilities or mixture parameters. For example, the special case in which each prior probability is fixed to the same value and all mixture components are assumed to have the same radius along all dimensions reduces to a soft version of the k -means algorithm. These connections will be discussed in detail in Chapter 3.

1.2.3 Distance-Based Algorithms

Many special forms of generative algorithms can be shown to reduce to distance-based algorithms. This is because the mixture components in generative models often use a distance function within the probability distribution. For example, the Gaussian distribution represents data generation probabilities in terms of the euclidian distance from the mean of the mixture. As a result, a generative model with the Gaussian distribution can be shown to have a very close relationship with the k -means algorithm. In fact, many distance-based algorithms can be shown to be reductions from or simplifications of different kinds of generative models.

Distance-based methods are often desirable because of their simplicity and ease of implementation in a wide variety of scenarios. Distance-based algorithms can be generally divided into two types:

- *Flat*: In this case, the data is divided into several clusters in one shot, typically with the use of partitioning representatives. The choice of the partitioning representative and distance function is crucial and regulates the behavior of the underlying algorithm. In each iteration, the data points are assigned to their closest partitioning representatives, and then the representative is adjusted according to the data points assigned to the cluster. It is instructive to compare this with the iterative nature of the EM algorithm, in which soft assignments are performed in the E-step, and model parameters (analogous to cluster representatives) are adjusted in the M-step. Some common methods for creating the partitions are as follows:

- *k-Means*: In these methods, the partitioning representatives correspond to the mean of each cluster. Note that the partitioning representative is not drawn from the original data set, but is created as a function of the underlying data. The euclidian distance is used in order to compute distances. The *k-Means* method is considered one of the simplest and most classical methods for data clustering [46] and is also perhaps one of the most widely used methods in practical implementations because of its simplicity.
- *k-Medians*: In these methods, the median along each dimension, instead of the mean, is used to create the partitioning representative. As in the case of the *k-Means* approach, the partitioning representatives are not drawn from the original data set. The *k-Medians* approach is more stable to noise and outliers, because the median of a set of values is usually less sensitive to extreme values in the data. It should also be noted that the term “*k-Medians*” is sometimes overloaded in the literature, since it is sometimes also used to refer to a *k-Medoid* approach (discussed below) in which the partitioning representatives are drawn from the original data. In spite of this overloading and confusion in the research literature, it should be noted that the *k-Medians* and *k-Medoid* methods should be considered as distinct techniques. Therefore, in several chapters of this book, the *k-Medians* approach discussed is really a *k-Medoids* approach, though we have chosen to be consistent within the specific research paper which is being described. Nevertheless, it would be useful to note the overloading of this term in order to avoid confusion.
- *k-Medoids*: In these methods, the partitioning representative is sampled from the original data. Such techniques are particularly useful in cases, where the data points to be clustered are arbitrary objects, and it is often not meaningful to talk about functions of these objects. For example, for a set of network or discrete sequence objects, it may not be meaningful to talk about their mean or median. In such cases, partitioning representatives are drawn from the data, and iterative methods are used in order to improve the quality of these representatives. In each iteration, one of the representatives is replaced with a representative from the current data, in order to check if the quality of the clustering improves. Thus, this approach can be viewed as a kind of hill climbing method. These methods generally require many more iterations than *k-Means* and *k-Medoids* methods. However, unlike the previous two methods, they can be used in scenarios where it is not meaningful to talk about means or medians of data objects (eg. structural data objects).
- *Hierarchical*: In these methods, the clusters are represented hierarchically through a *dendrogram*, at varying levels of granularity. Depending upon whether this hierarchical representation is created in top-down or bottom-up fashion, these representations may be considered either agglomerative or divisive.
 - *Agglomerative*: In these methods, a bottom-up approach is used, in which we start off with the individual data points and successively merge clusters in order to create a tree-like structure. A variety of choices are possible in terms of how these clusters may be merged, which provide different tradeoffs between quality and efficiency. Some examples of these choices are *single-linkage*, *all-pairs linkage*, *centroid-linkage*, and *sampled-linkage* clustering. In *single-linkage* clustering, the shortest distance between any pair of points in two clusters is used. In *all-pairs linkage*, the average over all pairs is used, whereas in *sampled linkage*, a sampling of the data points in the two clusters is used for calculating the average distance. In *centroid-linkage*, the distance between the centroids is used. Some variations of these methods have the disadvantage of *chaining*, in which larger clusters are naturally biased toward having closer distances to other points and will therefore attract a successively larger number of points. *Single linkage*

clustering is particularly susceptible to this phenomenon. A number of data domains such as network clustering are also more susceptible to this behavior.

- *Divisive*: In these methods, a top-down approach is used in order to successively partition the data points into a tree-like structure. Any flat clustering algorithm can be used in order to perform the partitioning at each step. Divisive partitioning allows greater flexibility in terms of both the hierarchical structure of the tree and the level of balance in the different clusters. It is not necessary to have a perfectly balanced tree in terms of the depths of the different nodes or a tree in which the degree of every branch is exactly two. This allows the construction of a tree structure which allows different tradeoffs in the balancing of the node depths and node weights (number of data points in the node). For example, in a top-down method, if the different branches of the tree are unbalanced in terms of node weights, then the largest cluster can be chosen preferentially for division at each level. Such an approach [53] is used in *METIS* in order to create well balanced clusters in large social networks, in which the problem of cluster imbalance is particularly severe. While *METIS* is not a distance-based algorithm, these general principles apply to distance-based algorithms as well.

Distance-based methods are very popular in the literature, because they can be used with almost any data type, as long as an appropriate distance function is created for that data type. Thus, the problem of clustering can be reduced to the problem of finding a distance function for that data type. Therefore, distance function design has itself become an important area of research for data mining in its own right [5, 82]. Dedicated methods also have often been designed for specific data domains such as categorical or time series data [32, 42]. Of course, in many domains, such as high-dimensional data, the quality of the distance functions reduces because of many irrelevant dimensions [43] and may show both errors and concentration effects, which reduce the statistical significance of data mining results. In such cases, one may use either the redundancy in larger portions of the pairwise distance matrix to abstract out the noise in the distance computations with spectral methods [19] or projections in order to directly find the clusters in relevant subsets of attributes [9]. A discussion of many hierarchical and partitioning algorithms is provided in Chapter 4 of this book.

1.2.4 Density- and Grid-Based Methods

Density- and grid-based methods are two closely related classes, which try to explore the data *space* at high levels of granularity. The density at any particular point in the data space is defined either in terms of the number of data points in a prespecified volume of its locality or in terms of a smoother kernel density estimate [74]. Typically, the data space is explored at a reasonably high level of granularity and a postprocessing phase is used in order to “put together” the dense regions of the data space into an arbitrary shape. Grid-based methods are a specific class of density-based methods in which the individual regions of the data space which are explored are formed into a grid-like structure. Grid-like structures are often particularly convenient because of greater ease in putting together the different dense blocks in the post-processing phase. Such grid-like methods can also be used in the context of high-dimensional methods, since the lower dimensional grids define clusters on subsets of dimensions [6].

A major advantage of these methods is that since they explore the data space at a high level of granularity, they can be used to reconstruct the entire shape of the data distribution. Two classical methods for density-based methods and grid-based methods are DBSCAN [34] and STING [83], respectively. The major challenge of density-based methods is that they are naturally defined on data points in a continuous space. Therefore, they often cannot be meaningfully used in a discrete or noneuclidian space, unless an embedding approach is used. Thus, many arbitrary data types such as time-series data are not quite as easy to use with density-based methods without specialized transformations. Another issue is that density computations becomes increasingly difficult to define

with greater dimensionality because of the greater number of cells in the underlying grid structure and the sparsity of the data in the underlying grid. A detailed discussion of density-based and grid-based methods is provided in Chapters 5 and 6 of this book.

1.2.5 Leveraging Dimensionality Reduction Methods

Dimensionality reduction methods are closely related to both feature selection and clustering, in that they attempt to use the closeness and correlations between dimensions to reduce the dimensionality of representation. Thus, dimensionality reduction methods can often be considered a vertical form of clustering, in which columns of the data are clustered with the use of either correlation or proximity analysis, as opposed to the rows. Therefore, a natural question arises, as to whether it is possible to perform these steps *simultaneously*, by clustering rows and columns of the data together. The idea is that simultaneous row and column clustering is likely to be more effective than performing either of these steps individually. This broader principle has led to numerous algorithms such as matrix factorization, spectral clustering, probabilistic latent semantic indexing, and co-clustering. Some of these methods such as spectral clustering are somewhat different but are nevertheless based on the same general concept. These methods are also closely related to projected clustering methods, which are commonly used for high dimensional data in the database literature. Some common models will be discussed below.

1.2.5.1 Generative Models for Dimensionality Reduction

In these models, a generative probability distribution is used to model the relationships between the data points and dimensions in an integrated way. For example, a generalized Gaussian distribution can be considered a mixture of arbitrarily correlated (oriented) clusters, whose parameters can be learned by the EM-algorithm. Of course, this is often not easy to do robustly with increasing dimensionality due to the larger number of parameters involved in the learning process. It is well known that methods such as EM are highly sensitive to overfitting, in which the number of parameters increases significantly. This is because EM methods try to retain *all* the information in terms of soft probabilities, rather than making the hard choices of point and dimension selection by non-parametric methods. Nevertheless, many special cases for different data types have been learned successfully with generative models.

A particularly common dimensionality reduction method is that of topic modeling in text data [45]. This method is also sometimes referred to as *Probabilistic Latent Semantic Indexing (PLSI)*. In topic modeling, a cluster is associated with a set of words and a set of documents simultaneously. The main parameters to be learned are the probabilities of assignments of words (dimensions) to topics (clusters) and those of the documents (data points) to topics (clusters). Thus, this naturally creates a soft clustering of the data from *both* a row and column perspective. These are then learned in an integrated way. These methods have found a lot of success in the text mining literature, and many methods, such as *Latent Dirichlet Allocation (LDA)*, which vary on this principle, with the use of more generalized priors have been proposed [22]. Many of these models are discussed briefly in Chapter 3.

1.2.5.2 Matrix Factorization and Co-Clustering

Matrix factorization and co-clustering methods are also commonly used classes of dimensionality reduction methods. These methods are usually applied to data which is represented as sparse nonnegative matrices, though it is possible in principle to generalize these methods to other kinds of matrices as well. However, the real attraction of this approach is the additional *interpretability* inherent in *nonnegative* matrix factorization methods, in which a data point can be expressed as a *nonnegative* linear combination of the concepts in the underlying data. Nonnegative matrix factor-

ization methods are closely related to co-clustering, which clusters the rows and columns of a matrix simultaneously.

Let A be a nonnegative $n \times d$ matrix, which contains n data entries, each of which has a dimensionality of d . In most typical applications such as text data, the matrix A represents small nonnegative quantities such as word frequencies and is not only nonnegative, but also sparse. Then, the matrix A can be *approximately* factorized into two nonnegative low rank matrices U and V of sizes $n \times k$ and $k \times d$, respectively. As we will discuss later, these matrices are the representatives for the clusters on the rows and the columns, respectively, when exactly k clusters are used in order to represent both rows and columns. Therefore, we have

$$A \approx U \cdot V \quad (1.1)$$

The residual matrix R represents the noise in the underlying data:

$$R = A - U \cdot V \quad (1.2)$$

Clearly, it is desirable to determine the factorized matrices U and V , such that the sum of the squares of the residuals in R is minimized. This is equivalent to determining nonnegative matrices U and V , such that the Froebinius norm of $A - U \cdot V$ is minimized. This is a constrained optimization problem, in which the constraints correspond to the nonnegativity of the entries in U and V . Therefore, a Lagrangian method can be used to learn the parameters of this optimization problem. A detailed discussion of the iterative approach is provided in [55].

The nonnegative $n \times k$ matrix U represents the coordinates of each of the n data points into each of the k newly created dimensions. A high positive value of the entry (i, j) implies that data point i is closely related to the newly created dimension j . Therefore, a trivial way to perform the clustering would be to assign each data point to the newly created dimension for which it has the largest component in U . Alternatively, if data points are allowed to belong to multiple clusters, then for each of the k columns in V , the entries with value above a particular threshold correspond to the document clusters. Thus, the newly created set of dimensions can be made to be synonymous with the clustering of the data set. In practice, it is possible to do much better, by using k -means on the new representation. One nice characteristic of the nonnegative matrix factorization method is that the size of an entry in the matrix U tells us how much a particular data point is related to a particular concept (or newly created dimension). The price of this greater interpretability is that the newly created sets of dimensions are typically not orthogonal to one another. This brings us to a discussion of the physical interpretation of matrix V .

The $k \times d$ matrix V provides the actual representation of each of the k newly created dimensions in terms of the original d dimensions. Thus, each row in this matrix is one component of this newly created axis system, and the rows are not necessarily orthogonal to one another (unlike most other dimensionality reduction methods). A large positive entry implies that this newly created dimension is highly related to a particular dimension in the underlying data. For example, in a document clustering application, the entries with large positive values in each row represent the most common words in each cluster. A thresholding technique can be used to identify these words. Note the similarity of this approach with a projection-based technique or the softer PLSI technique. In the context of a document clustering application, these large positive entries provide the *word clusters* in the underlying data. In the context of text data, each document can therefore be approximately expressed (because of the factorization process) as a nonnegative linear combination of at most k word-cluster vectors. The specific weight of that component represents the importance of that component, which makes the decomposition highly interpretable. Note that this interpretability is highly dependent on nonnegativity. Conversely, consider the word-membership vector across the corpus. This can be expressed in terms of at most k document-cluster vectors. This provides an idea of how important each document cluster is to that word.

At this point, it should be evident that the matrices U and V *simultaneously* provide the clusters

on the rows (documents) and columns (words). This general principle, when applied to sparse non-negative matrices, is also referred to as co-clustering. Of course, nonnegative matrix factorization is only one possible way to perform the co-clustering. A variety of graph-based spectral methods and other information theoretic methods can also be used in order to perform co-clustering [29, 30, 71]. Matrix factorization methods are discussed in Chapter 7. These techniques are also closely related to *spectral methods* for dimensionality reduction, as discussed below.

1.2.5.3 Spectral Methods

Spectral methods are an interesting technique for dimensionality reduction, which work with the similarity (or distance) matrix of the underlying data, instead of working with the original points and dimensions. This, of course, has its own set of advantages and disadvantages. The major advantage is that it is now possible to work with arbitrary objects for dimensionality reduction, rather than simply data points which are represented in a multi-dimensional space. In fact, spectral methods also perform the dual task of *embedding* these objects into a euclidian space, while performing the dimensionality reduction. Therefore, spectral methods are extremely popular for performing clustering on arbitrary objects such as node sets in a graph. The disadvantage of spectral methods is that since they work with an $n \times n$ similarity matrix, the time complexity for even creating the similarity matrix scales with the *square* of the number of *data points*. Furthermore, the process of determining the eigenvectors of this matrix can be extremely expensive, unless a very small number of these eigenvectors is required. Another disadvantage of spectral methods is that it is much more difficult to create lower dimensional representations for data points, unless they are part of the original sample from which the similarity matrix was created. For multidimensional data, the use of such a large similarity matrix is rather redundant, unless the data is extremely noisy and high dimensional.

Let \mathcal{D} be a database containing n points. The first step is to create an $n \times n$ matrix W of weights, which represents the pairwise similarity between the different data points. This is done with the use of the *heat kernel*. For any pair of data points \bar{X}_i and \bar{X}_j , the *heat kernel* defines a similarity matrix W of weights:

$$W_{ij} = \exp(-\|\bar{X}_i - \bar{X}_j\|^2/t) \quad (1.3)$$

Here t is a user-defined parameter. Furthermore, the value of W_{ij} is set to 0 if the distance between \bar{X}_i and \bar{X}_j is greater than a given threshold. The similarity matrix may also be viewed as the adjacency matrix of a graph, in which each node corresponds to a data item, and the weight of an edge corresponds to the similarity between these data items. Therefore, spectral methods reduce the problem to that of finding optimal cuts in this graph, which correspond to partitions which are weakly connected by edges representing similarity. Hence, spectral methods can be considered a graph-based technique for clustering of any kinds of data, by converting the similarity matrix into a network structure. Many variants exist in terms of the different choices for constructing the similarity matrix W . Some simpler variants use the mutual k -nearest neighbor graph, or simply the binary graph in which the distances are less than a given threshold. The matrix W is symmetric.

It should be noted that even when distances are very noisy, the similarity matrix encodes a significant amount of information because of its exhaustive nature. It is here that spectral analysis is useful, since the noise in the similarity representation can be abstracted out with the use of eigenvector-analysis of this matrix. Thus, these methods are able to recover and sharpen the latent information in the similarity matrix, though at a rather high cost, which scales with the square of the number of data points.

First, we will discuss the problem of mapping the points onto a 1-dimensional space. The generalization to the k -dimensional case is relatively straightforward. We would like to map the data points in \mathcal{D} into a set of points $y_1 \dots y_n$ on a line, in which the similarity maps onto euclidian distances on this line. Therefore, it is undesirable for data points which are very similar to be mapped

onto distant points on this line. We would like to determine values of y_i which minimize the following objective function O :

$$O = \sum_{i=1}^n \sum_{j=1}^n W_{ij} \cdot (y_i - y_j)^2 \quad (1.4)$$

The objective function O can be rewritten in terms of the Laplacian matrix L of W . The Laplacian matrix is defined as $D - W$, where D is a diagonal matrix satisfying $D_{ii} = \sum_{j=1}^n W_{ij}$. Let $\bar{y} = (y_1 \dots y_n)$. The objective function O can be rewritten as follows:

$$O = 2 \cdot \bar{y}^T \cdot L \cdot \bar{y} \quad (1.5)$$

We need to incorporate a scaling constraint in order to ensure that the trivial value of $y_i = 0$ for all i is not selected by the problem. A possible scaling constraint is as follows:

$$\bar{y}^T \cdot D \cdot \bar{y} = 1 \quad (1.6)$$

Note that the use of the matrix D provides different weights to the data items, because it is assumed that nodes with greater similarity to different data items are more involved in the clustering process. This optimization formulation is in generalized eigenvector format, and therefore the value of \bar{y} is optimized by selecting the smallest eigenvalue for which the generalized eigenvector relationship $L \cdot \bar{y} = \lambda \cdot D \cdot \bar{y}$ is satisfied. In practice however, the smallest generalized eigenvalue corresponds to the trivial solution, where \bar{y} is the (normalized) unit vector. This trivial eigenvector is noninformative. Therefore, it can be discarded, and it is not used in the analysis. The second-smallest eigenvalue then provides an optimal solution, which is more informative.

This model can be generalized to determining all the eigenvectors in ascending order of eigenvalue. Such directions correspond to successive orthogonal directions in the data, which result in the best mapping. This results in a set of n eigenvectors $\bar{e}_1, \bar{e}_2 \dots \bar{e}_n$ (of which the first is trivial), with corresponding eigenvalues $0 = \lambda_1 \leq \lambda_2 \dots \leq \lambda_n$. Let the corresponding vector representation of the data points along each eigenvector be denoted by $\bar{y}^1 \dots \bar{y}^n$.

What do the small and large magnitude eigenvectors intuitively represent in the new transformed space? By using the ordering of the data items along a small magnitude eigenvector to create a cut, the weight of the edges across the cut is likely to be small. Thus, this represents a cluster in the space of data items. At the same time, if the top k longest eigenvectors are picked, then the vector representations $\bar{y}^1 \dots \bar{y}^{n-k+1}$ provide a $n \times k$ matrix. This provides a k -dimensional embedded representation of the entire data set of n points, which preserves the maximum amount of information. Thus, spectral methods can be used in order to simultaneously perform clustering and dimensionality reduction. In fact, spectral methods can be used to recover the entire lower dimensional (possibly nonlinear) shape of the data, though arguably at a rather high cost [78]. The specific local dimensionality reduction technique of this section is discussed in detail in [19]. An excellent survey on spectral clustering methods may be found in [35, 63]. These methods are also discussed in detail in Chapter 8 of this book.

1.2.6 The High Dimensional Scenario

The high dimensional scenario is particularly challenging for cluster analysis because of the large variations in the behavior of the data attributes over different parts of the data. This leads to numerous challenges in many data mining problems such as clustering, nearest neighbor search, and outlier analysis. It should be noted that many of the algorithms for these problems are dependent upon the use of distances as an important subroutine. However, with increasing dimensionality, the distances seem to increasingly lose their effectiveness and statistical significance because of irrelevant attributes. The premise is that a successively smaller *fraction* of the attributes often remains relevant with increasing dimensionality, which leads to the blurring of the distances and increasing

concentration effects, because of the averaging behavior of the irrelevant attributes. *Concentration effects* refer to the fact that when many features are noisy and uncorrelated, their additive effects will lead to all pairwise distances between data points becoming similar. The noise and concentration effects are problematic in two ways for distance-based clustering (and many other) algorithms:

- An increasing noise from irrelevant attributes may cause errors in the distance representation, so that it no longer properly represents the *intrinsic distance* between data objects.
- The concentration effects from the irrelevant dimensions lead to a reduction in the statistical significance of the results from distance-based algorithms, if used directly with distances that have not been properly denoised.

The combination of these issues above leads to questions about whether full-dimensional distances are truly meaningful [8, 21, 43]. While the natural solution to such problems is to use feature selection, the problem is that different attributes may be relevant to different localities of the data. This problem is inherent in high-dimensional distance functions and nearest neighbor search. As stated in [43]: "... *One of the problems of the current notion of nearest neighbor search is that it tends to give equal treatment to all features (dimensions), which are however not of equal importance. Furthermore, the importance of a given dimension may not even be independent of the query point itself*" p. 506. These noise and concentration effects are therefore a problematic *symptom* of (locally) irrelevant, uncorrelated, or noisy attributes, which tend to impact the effectiveness and statistical significance of full-dimensional algorithms. This has led to significant efforts to redesign many data mining algorithms such as clustering, which are dependent on the notion of distances [4]. In particular, it would seem odd that data mining algorithms should behave poorly with increasing dimensionality at least from a qualitative perspective when a larger number of dimensions clearly provides more information. The reason is that conventional distance measures were generally designed for many kinds of low-dimensional spatial applications which are not suitable for the high-dimensional case. While high-dimensional data contain more information, they are also more complex. Therefore, naive methods will do worse with increasing dimensionality because of the noise effects of locally irrelevant attributes. Carefully designed distance functions can leverage the greater information in these dimensions and show *improving* behavior with dimensionality [4, 11] at least for a few applications such as similarity search.

Projected clustering methods can be considered a form of *local feature selection*, or *local dimensionality reduction*, in which the feature selection or transformation is performed specific to different localities of the data. Some of the earliest methods even refer to these methods as local-dimensionality reduction [23] in order to emphasize the local feature selection effect. Thus, a projected cluster is defined as a set of clusters $C_1 \dots C_k$, along with a corresponding set of subspaces $\mathcal{E}_1 \dots \mathcal{E}_k$, such that the points in C_i cluster well in the subspace represented by \mathcal{E}_i . Thus, these methods are a form of local feature selection, which can be used to determine the relevant clusters in the underlying data. Note that the subspace \mathcal{E}_i could represent a subset of the original attributes, or it could represent a transformed axis system in which the clusters are defined on a small set of orthogonal attributes. Some of the earliest projected clustering methods are discussed in [6, 9, 10]. The literature on this area has grown significantly since 2000, and surveys on the area may be found in [67, 54]. An overview of algorithms for high-dimensional data will be provided in Chapter 9.

It should also be pointed out that while pairwise distances are often too noisy or concentrated to be used meaningfully with off-the-shelf distance-based algorithms, larger portions of the *entire* distance matrix often retain a significant amount of latent structure due to the inherent redundancies in representing different pairwise distances. Therefore, even when there is significant noise and concentration of the distances, there will always be some level of consistent variations over the similarity matrix because of the impact of the consistent attributes. This redundancy can be leveraged in conjunction with spectral analysis [19, 78] to filter out the noise and concentration effects from full-dimensional distances and implicitly recover the local or global lower dimensional shape of the

underlying data in terms of enhanced distance representations of newly embedded data in a lower dimensional space. In essence, this approach enhances the contrasts of the distances by carefully examining how the noise (due to the many irrelevant attributes) and the minute correlations (due to the smaller number of relevant attributes) relate to the different pairwise distances. The general principle of these techniques is that distances are measured differently in high-dimensional space, depending upon how the data is distributed. This in turn depends upon the relevance and relationships of different attributes in different localities. These results are strong evidence for the fact that proper distance function design is highly dependent on its ability to recognize the relevance and relationships of different attributes in different data localities.

Thus, while (naively designed) pairwise distances cannot be used meaningfully in high dimensional data with *off-the-shelf* algorithms because of noise and concentration effects, they often do retain sufficient *latent* information *collectively* when used carefully in conjunction with denoising methods such as spectral analysis. Of course, the price for this is rather high, since the size of the similarity matrix scales with the square of the number of data points. The projected clustering method is generally a more efficient way of achieving an approximately similar goal, since it works directly with the data representation, rather than building a much larger and more redundant intermediate representation such as the distance matrix.

1.2.7 Scalable Techniques for Cluster Analysis

With advances in software and hardware technology, data collection has become increasingly easy in a wide variety of scenarios. For example, in social sensing applications, users may carry mobile or wearable sensors, which may result in the continuous accumulation of data over time. This leads to numerous challenges when real-time analysis and insights are required. This is referred to as the *streaming* scenario, in which it is assumed that a single pass is allowed over the data stream, because the data are often too large to be collected within limited resource constraints. Even when the data are collected offline, this leads to numerous scalability issues, in terms of integrating with traditional database systems or in terms of using the large amounts of data in a distributed setting, with the big data framework. Thus, varying levels of challenges are possible, depending upon the nature of the underlying data. Each of these issues is discussed below.

1.2.7.1 I/O Issues in Database Management

The most fundamental issues of scalability arise when a data mining algorithm is coupled with a traditional database system. In such cases, it can be shown that the major bottleneck arises from the I/O times required for accessing the objects in the database. Therefore, algorithms which use sequential scans of the data, rather than methods which randomly access the data records, are often useful. A number of classical methods were proposed in the database literature in order to address these scalability issues.

The easiest algorithms to extend to this case are flat partitioning methods which use sequential scans over the database in order to assign data points to representatives. One of the first methods [66] in this direction was *CLARANS*, in which these representatives are determined with the use of a k -medoids approach. Note that the k -medoids approach can still be computationally quite intensive because each iteration requires trying out new partitioning representatives through an exchange process (from the current set). If a large number of iterations is required, this will increase the number of passes over the data set. Therefore, the *CLARANS* method uses sampling, by performing the iterative hill climbing over a smaller sample, in order to improve the efficiency of the approach. Another method in this direction is *BIRCH* [87], which generalizes a k -means approach to the clustering process. The *CURE* method proposed in [41] finds clusters of nonspherical shape by using more than one representative point per cluster. It combines partitioning and sampling to ensure a

more efficient clustering process. A number of scalable algorithms for clustering are discussed in Chapter 11.

1.2.7.2 Streaming Algorithms

The streaming scenario is particularly challenging for clustering algorithms due to the requirements of *real-time* analysis, and the *evolution* and *concept-drift* in the underlying data. While database-centric algorithms require a *limited* number of passes over the data, streaming algorithms require exactly one pass, since the data cannot be stored at all. In *addition to* this challenge, the analysis typically needs to be performed in real time, and the changing patterns in the data need to be properly accounted for in the analysis.

In order to achieve these goals, virtually all streaming methods use a summarization technique to create *intermediate representations*, which can be used for clustering. One of the first methods in this direction uses a *microclustering approach* [7] to create and maintain the clusters from the underlying data stream. Summary statistics are maintained for the microclusters to enable an effective clustering process. This is combined with a pyramidal time frame to capture the evolving aspects of the underlying data stream. Stream clustering can also be extended to other data types such as discrete data, massive-domain data, text data, and graph data. A number of unique challenges also arises in the distributed setting. A variety of stream clustering algorithms is discussed in detail in Chapter 10.

1.2.7.3 The Big Data Framework

While streaming algorithms work under the assumption that the data are too large to be stored explicitly, the big data framework leverages advances in storage technology in order to actually store the data and process them. However, as the subsequent discussion will show, even if the data can be explicitly stored, it is often not easy to process and extract insights from them. This is because an increasing size of the data implies that a distributed file system must be used in order to store the information, and distributed processing techniques are required in order to ensure sufficient scalability. The challenge here is that if large segments of the data are available on different machines, it is often too expensive to shuffle the data across different machines in order to extract integrated insights from them. Thus, as in all distributed infrastructures, it is desirable to exchange intermediate insights, so as to minimize communication costs. For an application programmer, this can sometimes create challenges in terms of keeping track of where different parts of the data are stored, and the precise ordering of communications in order to minimize the costs.

In this context, Google's *MapReduce* framework [28] provides an effective method for analysis of large amounts of data, especially when the nature of the computations involves linearly computable statistical functions over the elements of the data streams. One desirable aspect of this framework is that it abstracts out the precise details of where different parts of the data are stored to the application programmer. As stated in [28]: "*The run-time system takes care of the details of partitioning the input data, scheduling the program's execution across a set of machines, handling machine failures, and managing the required inter-machine communication. This allows programmers without any experience with parallel and distributed systems to easily utilize the resources of a large distributed system*" p. 107. Many clustering algorithms such as *k*-means are naturally linear in terms of their scalability with the size of the data. A primer on the *MapReduce* framework implementation on *Apache Hadoop* may be found in [84]. The key idea here is to use a *Map* function to distribute the work across the different machines, and then provide an automated way to shuffle out much smaller data in (key,value) pairs containing intermediate results. The *Reduce* function is then applied to the aggregated results from the *Map* step to obtain the final results.

Google's original *MapReduce* framework was designed for analyzing large amounts of web logs and more specifically deriving linearly computable statistics from the logs. While the clustering process is not as simple as linearly computable statistics, it has nevertheless been shown [26] that many

existing clustering algorithms can be generalized to the *MapReduce* framework. A proper choice of the algorithm to adapt to the *MapReduce* framework is crucial, since the framework is particularly effective for linear computations. It should be pointed out that the major attraction of the *MapReduce* framework is its ability to provide application programmers with a cleaner abstraction, which is independent of very specific run-time details of the distributed system. It should not, however, be assumed that such a system is somehow inherently superior to existing methods for distributed parallelization from an *effectiveness* or *flexibility* perspective, especially if an application programmer is willing to design such details from scratch. A detailed discussion of clustering algorithms for big data is provided in Chapter 11.

1.3 Data Types Studied in Cluster Analysis

The specific data type has a tremendous impact on the particular choice of the clustering algorithm. Most of the earliest clustering algorithms were designed under the implicit assumption that the data attributes were numerical. However, this is not true in most real scenarios, where the data could be drawn from any number of possible types such as discrete (categorical), temporal, or structural. This section discusses the impact of data types on the clustering process. A brief overview of the different data types is provided in this section.

1.3.1 Clustering Categorical Data

Categorical data is fairly common in real data sets. This is because many attributes in real data such as sex, race, or zip code are inherently discrete and do not take on a natural ordering. In many cases, the data sets may be *mixed*, in which some attributes such as salary are numerical, whereas other attributes such as sex or zip code are categorical. A special form of categorical data is market basket data, in which all attributes are binary.

Categorical data sets lead to numerous challenges for clustering algorithms:

- When the algorithms depends upon the use of a similarity or distance function, the standard L_k metrics can no longer be used. New similarity measures need to be defined for categorical data. A discussion of similarity measures for categorical data is provided in [32].
- Many clustering algorithms such as the k -means or k -medians methods construct clustering representatives as the means or medians of the data points in the clusters. In many cases, statistics such as the mean or median are naturally defined for numerical data but need to be appropriately modified for discrete data.

When the data is mixed, then the problem because more difficult because the different attributes now need to be treated in a heterogeneous way, and the similarity functions need to explicitly account for the underlying heterogeneity.

It should be noted that some models of clustering are more amenable to different data types than others. For example, some models depend only on the distance (or similarity) functions between records. Therefore, as long as an appropriate similarity function can be defined between records, cluster analysis methods can be used effectively. Spectral clustering is one class of methods which can be used with virtually any data type, as long as appropriate similarity functions are defined. The downside is that the the methods scale with the square of the similarity matrix size. Generative models can also be generalized easily to different data types, as long as an appropriate generative model can be defined for each component of the mixture. Some common algorithms for categorical

data clustering include *CACTUS* [38], *ROCK* [40], *STIRR* [39], and *LIMBO* [15]. A discussion of categorical data clustering algorithms is provided in Chapter 12.

1.3.2 Clustering Text Data

Text data is a particularly common form of data with the increasing popularity of the web and social networks. Text data is typically represented in vector space format, in which the specific ordering is abstracted out, and the data is therefore treated as a bag-of-words. It should be noted that the methods for clustering text data can also be used for clustering set-based attributes. Text data has a number of properties which should be taken into consideration:

- The data is extremely high-dimensional and sparse. This corresponds to the fact that the text lexicon is rather large, but each document contains only a small number of words. Thus, most attributes take on zero values.
- The attribute values correspond to word frequencies and are, therefore, typically nonnegative. This is important from the perspective of many co-clustering and matrix factorization methods, which leverage this nonnegativity.

The earliest methods for text clustering such as the scatter-gather method [25, 75] use distance-based methods. Specifically, a combination of k -means and agglomerative methods is used for the clustering process. Subsequently, the problem of text clustering has often been explored in the context of topic modeling, where a soft membership matrix of words and documents to clusters is created. The EM-framework is used in conjunction with these methods. Two common methods used for generative topic modeling are *PLSI* and *LDA* [22, 45]. These methods can be considered soft versions of methods such as co-clustering [29, 30, 71] and matrix factorization [55], which cluster the rows and columns together at the same time. Spectral methods [73] are often used to perform this partitioning by creating a bipartite similarity graph, which represents the membership relations of the rows (documents) and columns (words). This is not particularly surprising since matrix factorization methods and spectral clustering are known to be closely related [57], as discussed in Chapter 8 of this book. Surveys on text clustering may be found in [12, 88]. In recent years, the popularity of social media has also lead to an increasing importance of *short* text documents. For example, the posts and tweets in social media web sites are constrained in length, both by the platform and by user preferences. Therefore, specialized methods have also been proposed recently for performing the clustering in cases where the documents are relatively short. Methods for clustering different kinds of text data are discussed in Chapter 13.

1.3.3 Clustering Multimedia Data

With the increasing popularity of social media, many forms of multimedia data may be used in conjunction with clustering methods. These include image data, audio and video. Examples of social media sites which contain large amounts of such data are *Flickr* and *Youtube*. Even the web and conventional social networks typically contain a significant amount of multimedia data of different types. In many cases, such data may occur in conjunction with other more conventional forms of text data.

The clustering of multimedia data is often a challenging task, because of the varying and heterogeneous nature of the underlying content. In many cases, the data may be multimodal, or it may be contextual, containing both behavioral and contextual attributes. For example, image data are typically contextual, in which the position of a pixel represents its context, and the value on the pixel represents its behavior. Video and music data are also contextual, because the temporal ordering of the data records provides the necessary information for understanding. The heterogeneity and contextual nature of the data can only be addressed with proper data representations and analysis. In

fact, data representation seems to be a key issue in all forms of multimedia analysis, which significantly impacts the final quality of results. A discussion of methods for clustering multimedia data is provided in Chapter 14.

1.3.4 Clustering Time-Series Data

Time-series data is quite common in all forms of sensor data, stock markets, or any other kinds of temporal tracking or forecasting applications. The major aspect of time series is that the data values are not independent of one another, but they are temporally dependent on one another. Specifically, the data contain a contextual attribute (time) and a behavioral attribute (data value). There is a significant diversity in problem definitions in the time-series scenario. The time-series data can be clustered in a variety of different ways, depending upon whether correlation-based online analysis is required or shape-based offline analysis is required.

- In correlation-based online analysis, the correlations among the different time-series data streams are tracked over time in order to create online clusters. Such methods are often useful for sensor selection and forecasting, especially when streams exhibit lag correlations within the clustered patterns. These methods are often used in stock market applications, where it is desirable to maintain groups of clustered stock tickers, in an online manner, based on their correlation patterns. Thus, the distance functions between different series need to be computed continuously, based on their mutual regression coefficients. Some examples of these methods include the MUSCLES approach [86] and a large scale time-series correlation monitoring method [90].
- In shape-based offline analysis, the time-series objects are analyzed in offline manner, and the specific details about when a particular time series was created is not important. For example, for a set of ECG time series collected from patients, the precise time of when a series was collected is not important, but the overall shape of the series is important for the purposes of clustering. In such cases, the distance function between two time series is important. This is important, because the different time series may not be drawn on the same range of data values and may also show time-warping effects, in which the shapes can be matched only by elongating or shrinking portions of the time-series in the temporal direction. As in the previous case, the design of the distance function [42] holds the key to the effective use of the approach.

A particular interesting case is that of multivariate time series, in which many series are simultaneously produced over time. A classical example of this is trajectory data, in which the different coordinate directions form the different components of the multivariate series. Therefore, trajectory analysis can be viewed as a special kind of time-series clustering. As in the case of univariate time series, it is possible to perform these steps using either online analysis (trajectories moving together in real time) or offline analysis (similar shape). An example of the former is discussed in [59], whereas an example of the latter is discussed in [68]. A survey on time-series data is found in [60], though this survey is largely focussed on the case of offline analysis. Chapter 15 discusses both the online and offline aspects of time series clustering.

1.3.5 Clustering Discrete Sequences

Many forms of data create discrete sequences instead of categorical ones. For example, web logs, the command sequences in computer systems, and biological data are all discrete sequences. The contextual attribute in this case often corresponds to placement (e.g., biological data), rather than time. Biological data is also one of the most common applications of sequence clustering.

As with the case of continuous sequences, a key challenge is the creation of similarity functions between different data objects. Numerous similarity functions such as the hamming distance, edit distance, and longest common subsequence are commonly used in this context. A discussion of the similarity functions which are commonly used for discrete sequences may be found in [58]. Another key problem which arises in the context of clustering discrete sequences is that the intermediate and summary representation of a set of sequences can be computationally intensive. Unlike numerical data, where averaging methods can be used, it is much more difficult to find such representations for discrete sequences. A common representation, which provides a relatively limited level of summarization is the suffix tree. Methods for using suffix trees for sequence clustering methods have been proposed in CLUSEQ [85].

Generative models can be utilized, both to model the distances between sequences and to create probabilistic models of cluster generation [76]. A particularly common approach is the use of mixtures of HMMs. A primer on Hidden Markov Models may be found in [70]. Hidden Markov Models can be considered a special kind of mixture model in which the different components of the mixture are dependent on one another. A second level of mixture modeling can be used in order to create clusters from these different HMMs. Much of the work on sequence clustering is performed in the context of biological data. Detailed surveys on the subject may be found in [33, 49, 64]. A discussion of sequence clustering algorithms is provided in Chapter 16, though this chapter also provides a survey of clustering algorithms for other kinds of biological data.

1.3.6 Clustering Network Data

Graphs and networks are among the most fundamental (and general) of all data representations. This is because virtually every data type can be represented as a similarity graph, with similarity values on the edges. In fact, the method of spectral clustering can be viewed as the most general connection between all other types of clustering and graph clustering. Thus, as long as a similarity function can be defined between arbitrary data objects, spectral clustering can be used in order to perform the analysis. Graph clustering has been studied extensively in the classical literature, especially in the context of the problem of 2-way and multi-way graph partitioning. The most classical of all multi-way partitioning algorithms is the Kernighan-Lin method [56]. Such methods can be used in conjunction with graph coarsening methods in order to provide efficient solutions. These methods are known as multilevel graph partitioning techniques. A particularly popular algorithm in this category is METIS [53].

A number of methods are commonly used in the literature in order to create partitions from graphs:

- *Generative Models:* As discussed earlier in this chapter, it is possible to define a generative model for practically any clustering problem, as long as an appropriate method exists for defining each component of the mixture as a probability distribution. An example of a generative model for network clustering is found in [77].
- *Classical Combinatorial Algorithms:* These methods use network flow [13] or other iterative combinatorial techniques in order to create partitions from the underlying graph. It should be pointed out that even edge sampling is often known to create good quality partitions, when it is repeated many times [51]. It is often desirable to determine cuts which are well balanced across different partitions, because the cut with the smallest absolute value often contains the large majority of the nodes in a single partition and a very small number of nodes in the remaining partitions. Different kinds of objective functions in terms of creating the cut, such as the unnormalized cut, normalized cut, and ratio cut, provide different tradeoffs between cluster balance and solution quality [73]. It should be pointed out that since graph cuts are a combinatorial optimization problem, they can be formulated as integer programs.

- *Spectral Methods*: Spectral methods can be viewed as *linear programming relaxations* to the integer programs representing the optimization of graph cuts. Different objective functions can be constructed for different kinds of cuts, such as the unnormalized cut, ratio cut, and normalized cut. Thus, the continuous solutions to these linear programs can be used to create a multidimensional embedding for the nodes, on which conventional k -means algorithms can be applied. These linear programs can be shown to take on a specially convenient form, in which the generalized eigenvectors of the graph Laplacian correspond to solutions of the optimization problem.
- *Nonnegative Matrix Factorization*: Since a graph can be represented as an adjacency matrix, nonnegative matrix factorization can be used in order to decompose it into two low rank matrices. It is possible to apply the matrix factorization methods either to the node–node adjacency matrix or the node–edge adjacency matrix to obtain different kinds of insights. It is also relatively easy to augment the matrix with content to create analytical methods, which can cluster with a combination of content and structure [69].

While the aforementioned methods represent a sampling of the important graph clustering methods, numerous other objective functions are possible for the construction of graph cuts such as the use of modularity-based objective functions [24]. Furthermore, the problem becomes even more challenging in the context of social networks, where content may be available at either the nodes [89] or edges [69]. Surveys on network clustering may be found in [36, 72]. Algorithms for network clustering are discussed in detail in Chapter 17.

1.3.7 Clustering Uncertain Data

Many forms of data either are of low fidelity or the quality of the data has been intentionally degraded in order to design different kinds of network mining algorithms. This has led to the field of probabilistic databases. Probabilistic data can be represented either in the form of attribute-wise uncertainty or in the form of a *possible worlds* model, in which only particular subsets of attributes can be present in the data at a given time [3]. The key idea here is that the incorporation of probabilistic information can improve the quality of data mining algorithms. For example, if two attributes are equally desirable to use in an algorithm in the deterministic scenario, but one of them has greater uncertainty than the other, then the attribute with less uncertainty is clearly more desirable for use. Uncertain clustering algorithms have also been extended recently to the domain of streams and graphs. In the context of graphs, it is often desirable to determine the most reliable subgraphs in the underlying network. These are the subgraphs which are the most difficult to disconnect under edge uncertainty. A discussion of algorithms for reliable graph clustering may be found in [62]. Uncertain data clustering algorithms are discussed in Chapter 18.

1.4 Insights Gained from Different Variations of Cluster Analysis

While the aforementioned methods discuss the basics of the different clustering algorithms, it is often possible to obtain enhanced insights either by using more rigorous analysis or by incorporating additional techniques or data inputs. These methods are particularly important because of the subjectivity of the clustering process, and the many different ways in which the same data set can be clustered. How do we know that a particular clustering is good or that it solves the needs of the application? There are numerous ways in which these issues can be explored. The exploration could be through interactive visualization and human interaction, external knowledge-based supervision,

explicitly exploring the multiple solutions to evaluate different possibilities, combining the multiple solutions to create more robust ensembles, or trying to judge the quality of different solutions with the use of different validation criteria. For example, the presence of labels adds domain knowledge to the clustering process, which can be used in order to improve insights about the quality of the clustering. This approach is referred to as semisupervision. A different way of incorporating domain knowledge (and indirect supervision) would be for the human to explore the clusters interactively with the use of visual methods and interact with the clustering software in order to create more meaningful clusters. The following subsections will discuss these alternatives in detail.

1.4.1 Visual Insights

Visual analysis of multidimensional data is a very intuitive way to explore the structure of the underlying data, possibly incorporating human feedback into the process. Thus, this approach could be considered an informal type of supervision, when human feedback is incorporated into cluster generation. The major advantage of incorporating human interaction is that a human can often provide intuitive insights, which are not possible from an automated computer program of significant complexity. On the other hand, a computer is much faster at the detailed calculations required both for clustering and for “guessing” the most appropriate feature-specific views of high dimensional data. Thus, a combination of a human and a computer often provides clusters which are superior to those created by either.

One of the most well-known systems for visualization of high-dimensional clusters is the HD-Eye method [44], which explores different subspaces of the data in order to determine clusters in different feature-specific views of the data. Another well-known technique is the *IPCLUS* method [2]. The latter method generates feature-specific views in which the data is *well polarized*. A well-polarized view refers to a 2-dimensional subset of features in which the data clearly separates out into clusters. A kernel-density estimation method is used to determine the views in which the data is well polarized. The final clustering is determined by exploring different views of the data, and counting how the data separates out into clusters in these different views. This process is essentially an ensemble-based method, an approach which is used popularly in the clustering literature, and will be discussed in a later part of this section. Methods for both incorporating and extracting visual insights from the clustering process are discussed in Chapter 19.

1.4.2 Supervised Insights

The same data set may often be clustered in multiple ways, especially when the dimensionality of the data set is high and subspace methods are used. Different features may be more relevant to different kinds of applications and insights. Since clustering is often used as an intermediate step in many data mining algorithms, it then becomes difficult to choose a particular kind of clustering that may suit that application. The subjectiveness of clustering is highly recognized, and small changes in the underlying algorithm or data set may lead to significant changes in the underlying clusters. In many cases, this subjectiveness also implies that it is difficult to refer to one particular clustering as significantly better than another. It is here that supervision can often play an effective role, because it takes the specific goal of the analyst into consideration.

Consider a document clustering application, in which a web portal creator (analyst) wishes to segment the documents into a number of categories. In such cases, the analyst may already have an approximate idea of the categories in which he is interested, but he may not have *fully* settled on a particular set of categories. This is because the data may also contain as yet unknown categories in which the analyst is interested. In such cases, semisupervision is an appropriate way to approach the problem. A number of labeled examples are provided, which approximately represent the categories in which the analyst is interested. This is used as domain knowledge or prior knowledge,

which is used in order to supervise the clustering process. For example, a very simple form of supervision would be to use seeding, in which the documents of the appropriate categories are provided as (some of the) seeds to a representative clustering algorithm such as k -means. In recent years, spectral methods have also been heavily adapted for the problem of semisupervised clustering. In these methods, Laplacian smoothing is used on the labels to generate the clusters. This allows the learning of the lower dimensional data surface in a semisupervised way, as it relates to the underlying clusters. The area of semisupervised clustering is also sometimes referred to as constrained clustering. An excellent discussion on constrained clustering algorithms may be found in [18]. A number of interesting methods for semisupervised clustering are discussed in Chapter 20, with a special focus on the graph-based algorithms.

1.4.3 Multiview and Ensemble-Based Insights

As discussed above, one of the major issues in the clustering process is that different kinds of clusters are possible. When no supervision is available, the bewildering number of possibilities in the clustering process can sometimes be problematic for the analyst, especially from the perspective of interpretability. These are referred to as alternative clusters, and technically represent the behavior from different perspectives. In many cases, the ability to provide different clustering solutions that are significantly different provides insights to the analyst about the key clustering properties of the underlying data. This broader area is also referred to as multiview clustering.

The most naive method for multiview clustering is to simply run the clustering algorithm multiple times, and then examine the different clustering solutions to determine those which are different. A somewhat different approach is to use spectral methods in order to create approximately orthogonal clusters. Recall that the eigenvectors of the Laplacian matrix represent alternative cuts in the graph and that the small eigenvectors represent the best cuts. Thus, by applying a 2-means algorithm to the embedding on each eigenvector, it is possible to create a clustering which is very different from the clustering created by other (orthogonal) eigenvectors. The orthogonality of the eigenvectors is important, because it implies that the embedded representations are very different. Furthermore, the smallest eigenvectors represent the best clusterings, whereas clusterings derived from successively larger eigenvectors represent successively suboptimal solutions. Thus, this approach not only provides alternative clusterings which are quite different from one another, but also provides a ranking of the quality of the different solutions. A discussion of alternative clustering methods is provided in Chapter 21.

In many cases, the alternative clustering methods can be combined to create more robust solutions with the use of ensemble-based techniques. The idea here is that a combination of the output of the different clusterings provides a more robust picture of how the points are related to one another. Therefore, the outputs of the different alternatives can be used as input to a meta-algorithm which combines the results from the different algorithms. Such an approach provides a more robust clustering solution. A discussion of ensemble-based methods for clustering is provided in Chapter 22.

1.4.4 Validation-Based Insights

Given a particular clustering, how do we know what the quality of the clustering really is? While one possible approach is to use synthetic data to determine the matching between the input and output clusters, it is not fully satisfying to rely on only synthetic data. This is because the results on synthetic data may often be specific to a particular algorithm and may not be easily generalizable to arbitrary data sets.

Therefore, it is desirable to use validation criteria on the basis of real data sets. The problem in the context of the clustering problem is that the criteria for quality is not quite as crisp as many other data mining problems such as classification, where external validation criteria are available in

the form of labels. Therefore, the use of one or more criteria may inadvertently favor different algorithms. As the following discussion suggests, clustering is a problem in which precise quantification is often not possible because of its unsupervised nature. Nevertheless, many techniques provide a partial understanding of the underlying clusters. Some common techniques for cluster validation are as follows:

- A common method in the literature is to use case studies to illustrate the subjective quality of the clusters. While case studies provide good intuitive insights, they are not particularly effective for providing a more rigorous quantification of the quality. It is often difficult to compare two clustering methods from a quantitative perspective with the use of such an approach.
- Specific measurements of the clusters such as the cluster radius or density may be used in order to provide a measure of quality. The problem here is that these measures may favor different algorithms in a different way. For example, a k -means approach will typically be superior to a density-based clustering method in terms of average cluster radius, but a density-based method may be superior to a k -means algorithms in terms of the estimated density of the clusters. This is because there is a circularity in using a particular criterion to evaluate the algorithm, when the same criterion is used for clustering purposes. This results in a bias during the evaluation. On the other hand, it may sometimes be possible to reasonably compare two different algorithms of a very similar type (e.g., two variations of k -means) on the basis of a particular criterion.
- In many data sets, labels may be associated with the data points. In such cases, cluster quality can be measured in terms of the correlations of the clusters with the data labels. This provides an external validation criterion, if the labels have not been used in the clustering process. However, such an approach is not perfect, because the class labels may not always align with the natural clusters in the data. Nevertheless, the approach is still considered more “impartial” than the other two methods discussed above and is commonly used for cluster evaluation.

A detailed discussion of the validation methods commonly used in clustering algorithms is provided in Chapter 23.

1.5 Discussion and Conclusions

Clustering is one of the most fundamental data mining problems because of its numerous applications to customer segmentation, target marketing, and data summarization. Numerous classes of methods have been proposed in the literature, such as probabilistic methods, distance-based methods, density-based methods, grid-based methods, factorization techniques, and spectral methods. The problem of clustering has close relationships to the problems of dimensionality reduction, especially through projected clustering formulations. High-dimensional data is often challenging for analysis, because of the increasing sparsity of the data. Clustering methods can be viewed as an integration of feature selection/dimensionality reduction methods with clustering.

The increasing advances in hardware technology allow the collection of large amounts of data through in an ever-increasing number of ways. This requires dedicated methods for streaming and distributed processing. Streaming methods typically work with only one pass over the data, and explicitly account for temporal locality in the clustering methods. Big data methods develop distributed techniques for clustering, especially through the *MapReduce* framework. The diversity of different data types significantly adds to the richness of the clustering problems. Many variations and enhancements of clustering such as visual methods, ensemble methods, multiview methods, or

supervised methods can be used to improve the quality of the insights obtained from the clustering process.

Bibliography

- [1] C. Aggarwal. *Data Streams: Models and Algorithms*, Springer, 2007.
- [2] C. C. Aggarwal. A human-computer interactive system for projected clustering, *ACM KDD Conference*, 2001.
- [3] C. Aggarwal. *Managing and Mining Uncertain Data*, Springer, 2009.
- [4] C. Aggarwal. Re-designing distance functions and distance-based applications for high dimensional data, *ACM SIGMOD Record*, March 2001.
- [5] C. Aggarwal. Towards systematic design of distance functions for data mining applications, *KDD Conference*, 2003.
- [6] R. Agrawal, J. Gehrke, D. Gunopulos, and P. Raghavan. Automatic subspace clustering of high dimensional data for data mining applications. *ACM SIGMOD Conference*, 1998.
- [7] C. Aggarwal, J. Han, J. Wang, and P. Yu. A framework for clustering evolving data streams. In *VLDB Conference*, 2003.
- [8] C. Aggarwal, A. Hinneburg, and D. Keim. On the surprising behavior of distance metrics in high dimensional space, *ICDT Conference*, 2001.
- [9] C. Aggarwal, C. Procopiuc, J. Wolf, P. Yu, and J.-S. Park. Fast algorithms for projected clustering. In *ACM SIGMOD Conference*, 1999.
- [10] C. Aggarwal and P. Yu. Finding generalized projected clusters in high dimensional spaces, *ACM SIGMOD Conference*, 2000.
- [11] C. Aggarwal and P. Yu. The Igrid index: Reversing the dimensionality curse for similarity indexing in high dimensional space, *KDD Conference*, 2000.
- [12] C. Aggarwal and C. Zhai. A survey of text clustering algorithms, *Mining Text Data*, Springer, 2012.
- [13] R. Ahuja, T. Magnanti, and J. Orlin. *Network Flows: Theory, Algorithms and Applications*, Prentice Hall, 1993.
- [14] M. Anderberg. *Cluster Analysis for Applications*. Academic Press, 1973.
- [15] P. Andritsos et al. LIMBO: Scalable clustering of categorical data. *EDBT Conference*, 2004.
- [16] A. Baraldi and P. Blonda. A survey of fuzzy clustering algorithms for pattern recognition. I. *IEEE Transactions on Systems, Man, and Cybernetics, Part B: Cybernetics*, 29(6):778–785, 1999.
- [17] A. Baraldi and P. Blonda. A survey of fuzzy clustering algorithms for pattern recognition. II. *IEEE Transactions on Systems, Man, and Cybernetics, Part B: Cybernetics*, 29(6):786–801, 1999.

- [18] S. Basu, I. Davidson, and K. Wagstaff. *Constrained clustering: Advances in theory, Algorithms and applications*, Chapman and Hall, 2008.
- [19] M. Belkin and P. Niyogi. Laplacian eigenmaps for dimensionality reduction and data representation, *Neural Computation*, 15(6), 2003.
- [20] P. Berkhin. A survey of clustering data mining techniques. *Grouping Multidimensional Data*, pages 25–71, 2006.
- [21] K. Beyer, J. Goldstein, U. Shaft, and R. Ramakrishnan. When is nearest neighbor meaningful? *ICDT Conference*, 2001.
- [22] D. Blei, A. Ng, and M. Jordan. Latent Dirichlet allocation, *Journal of Machine Learning Research*, 3:993–1022, 2003.
- [23] K. Chakrabarti and S. Mehrotra. Local dimensionality reduction: A new approach to indexing high dimensional spaces. *VLDB Conference Proceedings*, pages 89–100, 2000.
- [24] A. Clauset, M. E. J. Newman, and C. Moore, Finding community structure in very large networks, *Physical Review E* 70:066111, 2004.
- [25] D. Cutting, D. Karger, J. Pedersen, and J. Tukey. Scatter/gather: A cluster-based approach to browsing large document collections. *ACM SIGIR Conference*, pages 318–329, 1992.
- [26] R. Cordeiro, C. Traina Jr., A. Traina, J. López, U. Kang, and C. Faloutsos. Clustering very large multi-dimensional datasets with mapreduce. In *KDD*, pages 690–698, 2011.
- [27] A. P. Dempster, N. M. Laird, and D. B. Rubin. Maximum likelihood from incomplete data via the EM algorithm, *Journal of the Royal Statistical Society*, B, 39(1):1–38, 1977.
- [28] J. Dean and S. Ghemawat. MapReduce: A flexible data processing tool, *Communication of the ACM*, 53: 72–77, 2010.
- [29] I. Dhillon. Co-clustering documents and words using bipartite spectral graph partitioning, *ACM KDD Conference*, 2001.
- [30] I. Dhillon, S. Mallela, and D. Modha. Information-theoretic co-clustering, *ACM KDD Conference*, 2003.
- [31] B. Duran. *Cluster Analysis: A Survey*. Springer-Verlag, 1974.
- [32] G. Das and H. Mannila. Context-based similarity measures for categorical databases. *PKDD Conference*, pages 201–210, 2000.
- [33] M. B. Eisen, P. T. Spellman, P. O. Brown, and D. Botstein. Cluster analysis and display of genome-wide expression patterns. *Proceedings of the National Academy of Sciences*, 95(25):14863–14868, 1998. <http://rana.lbl.gov/EisenSoftware.htm>
- [34] M. Ester, H.-P. Kriegel, J. Sander, and X. Xu. A density-based algorithm for discovering clusters in large spatial databases with noise. *ACM KDD Conference*, pages 226–231, 1996.
- [35] M. Filippone, F. Camastra, F. Masulli, and S. Rovetta. A survey of kernel and spectral methods for clustering. *Pattern Recognition*, 41(1):176–190, 2008.
- [36] S. Fortunato. Community detection in graphs, *Physics Reports*, 486(3–5):75–174, February 2010.

- [37] G. Gan, C. Ma, and J. Wu. *Data Clustering: Theory, Algorithms, and Applications*. SIAM, Society for Industrial and Applied Mathematics, 2007.
- [38] V. Ganti, J. Gehrke, and R. Ramakrishnan. CACTUS-clustering categorical data using summaries. *ACM KDD Conference*, 1999.
- [39] D. Gibson, J. Kleiberg, and P. Raghavan. Clustering categorical data: An approach based on Dynamical Systems. *VLDB Conference*, 1998.
- [40] S. Guha, R. Rastogi, and K. Shim. ROCK: A robust clustering algorithm for categorical attributes. *Information Systems*, 25(5):345–366, 2000.
- [41] S. Guha, R. Rastogi, and K. Shim. CURE: An efficient clustering algorithm for large databases. In *ACM SIGMOD Conference*, 1998.
- [42] D. Gunopulos and G. Das. Time-series similarity measures, and time series indexing, *ACM SIGMOD Conference*, 2001.
- [43] A. Hinneburg, C. Aggarwal, and D. Keim. What is the nearest neighbor in high dimensional spaces, *VLDB Conference*, 2000.
- [44] A. Hinneburg, D. Keim, and M. Wawryniuk. Hd-eye: Visual mining of high-dimensional data. *IEEE Computer Graphics and Applications*, 19:22–31, 1999.
- [45] T. Hofmann. Probabilistic latent semantic indexing. *ACM SIGIR Conference*, 1999.
- [46] A. Jain. Data clustering: 50 years beyond k -means. *Pattern Recognition Letters*, 31(8):651–666, 2010.
- [47] A. Jain and R. Dubes. *Algorithms for Clustering Data*. Prentice Hall, 1988.
- [48] A. Jain, M. Murty, and P. Flynn. Data clustering: A review. *ACM Computing Surveys (CSUR)*, 31(3):264–323, 1999.
- [49] D. Jiang, C. Tang, and A. Zhang. Cluster analysis for gene expression data: A survey. *IEEE Transactions on Knowledge and Data Engineering*, 16(11):1370–1386, 2004.
- [50] I. Jolliffe. *Principal Component Analysis*, Springer, 2002.
- [51] D. R. Karger. Random sampling in cut, flow, and network design problems, *STOC*, pp. 648–657, 1994.
- [52] L. Kaufman and P. Rousseeuw. *Finding Groups in Data: An Introduction to Cluster Analysis*. Wiley-Interscience, 2005.
- [53] G. Karypis and V. Kumar. Multilevel k -way partitioning scheme for irregular graphs. *Journal of Parallel and Distributed Computing*, 48(1):96–129, 1998.
- [54] H.-P. Kriegel, P. Kröger, and A. Zimek. Clustering high-dimensional data: A survey on subspace clustering, pattern-based clustering, and correlation clustering. *ACM Transactions on Knowledge Discovery from Data (TKDD)*, 3(1):1–58, 2009.
- [55] D. D. Lee and H. S. Seung. Learning the parts of objects by non-negative matrix factorization, *Nature*, 401:788–791, 1999.
- [56] B.W. Kernighan and S. Lin. An efficient heuristic procedure for partitioning graphs, *Bell System Tech. Journal*, 49:291–307, Feb. 1970.

- [57] C. Ding, X. He, and H. Simon. On the equivalence of nonnegative matrix factorization and spectral clustering. *SDM Conference*, 2005.
- [58] D. Gusfield. *Algorithms for Strings, Trees and Sequences*, Cambridge University Press, 1997.
- [59] J.-G. Lee, J. Han, and K.-Y. Whang. Trajectory clustering: A partition-and-group framework. *SIGMOD Conference*, 593–604, 2007.
- [60] T. Liao. Clustering of time series data—A survey. *Pattern Recognition*, 38(11):1857–1874, 2005.
- [61] H. Liu and H. Motoda. *Computational Methods of Feature Selection*, Chapman and Hall (CRC Press), 2007.
- [62] L. Liu, R. Jin, C. Aggarwal, and Y. Shen. Reliable clustering on uncertain graphs, *ICDM Conference*, 2012.
- [63] U. von Luxberg. A tutorial on spectral clustering, *Statistics and Computing*, Springer, 2007.
- [64] S. Madeira and A. Oliveira. Biclustering algorithms for biological data analysis: A survey. *IEEE/ACM Transactions on Computational Biology and Bioinformatics*, 1(1):24–45, 2004.
- [65] F. Murtagh. A survey of recent advances in hierarchical clustering algorithms. *The Computer Journal*, 26(4):354–359, 1983.
- [66] R. T. Ng and J. Han. CLARANS: A method for clustering objects for spatial data mining. *IEEE Trans. Knowl. Data Eng.*, 14(5):1003–1016, 2002
- [67] L. Parsons, E. Haque, and H. Liu. Subspace clustering for high dimensional data: A review. *ACM SIGKDD Explorations*, 6(1):90–105, 2004.
- [68] G. Qi, C. Aggarwal, and T. Huang. Online community detection in social sensing. *WSDM Conference*, 2013.
- [69] G. Qi, C. Aggarwal, and T. Huang. Community detection with edge content in social media networks, *ICDE Conference*, 2013.
- [70] L. R. Rabiner. A tutorial on hidden Markov models and selected applications in speech recognition, *Proceedings of the IEEE*, 77(2):257–285, Feb. 1989.
- [71] M. Rege, M. Dong, and F. Fotouhi. Co-clustering documents and words using bipartite isoperimetric graph partitioning. *ICDM Conference*, pp. 532–541, 2006.
- [72] S. Schaeffer. Graph clustering. *Computer Science Review*, 1(1):27–64, 2007.
- [73] J. Shi and J. Malik. Normalized cuts and image segmentation. *IEEE Transaction on Pattern Analysis and Machine Intelligence*, 2000.
- [74] B. W. Silverman. *Density Estimation for Statistics and Data Analysis*. Chapman and Hall, 1986.
- [75] H. Schutze and C. Silverstein. Projections for efficient document clustering, *ACM SIGIR Conference*, 1997.
- [76] P. Smyth. Clustering sequences with hidden Markov models, *Neural Information Processing*, 1997.

- [77] Y. Sun, C. Aggarwal, and J. Han. Relation-strength aware clustering of heterogeneous information networks with incomplete attributes, *Journal of Proceedings of the VLDB Endowment*, 5(5):394–405, 2012.
- [78] J. Tenenbaum, V. de Silva, and J. Langford. A global geometric framework for nonlinear dimensionality reduction. *Science*, 290(5500):2319–2323, 2000.
- [79] P. Willett. Recent trends in hierarchical document clustering: A critical review. *Information Processing & Management*, 24(5):577–597, 1988.
- [80] R. Xu and D. Wunsch. Survey of clustering algorithms. *IEEE Transactions on Neural Networks*, 16(3):645–678, 2005.
- [81] R. Xu and D. Wunsch. *Clustering*. Wiley-IEEE Press, 2008.
- [82] F. Wang and J. Sun. Distance metric learning in data mining, *SDM Conference (Tutorial)*, 2012.
- [83] W. Wang, J. Yang, and R. Muntz. Sting: A statistical information grid approach to spatial data mining. *VLDB Conference*, 1997.
- [84] T. White. Hadoop: The Definitive Guide. *Yahoo! Press*, 2011.
- [85] J. Yang and W. Wang. CLUSEQ: Efficient and effective sequence clustering, *ICDE Conference*, 2003.
- [86] B.-K. Yi, N. D. Sidiropoulos, T. Johnson, H. Jagadish, C. Faloutsos, and A. Biliris. Online data mining for co-evolving time sequences. *ICDE Conference*, 2000.
- [87] T. Zhang, R. Ramakrishnan, and M. Livny. BIRCH: An efficient data clustering method for very large databases. In *ACM SIGMOD Conference*, 1996.
- [88] Y. Zhao, G. Karypis, and U. Fayyad. Hierarchical clustering algorithms for document datasets. *Data Mining and Knowledge Discovery*, 10(2):141–168, 2005.
- [89] Y. Zhou, H. Cheng, and J. X. Yu. Graph clustering based on structural/attribute similarities, *Proc. VLDB Endow.*, 2(1):718–729, 2009.
- [90] Y. Zhu and D. Shasha. StatStream: Statistical monitoring of thousands of data streams in real time. *VLDB Conference*, pages 358–369, 2002.