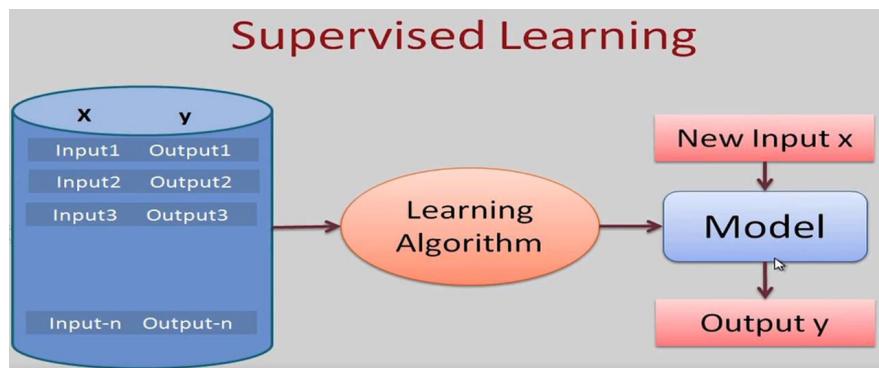
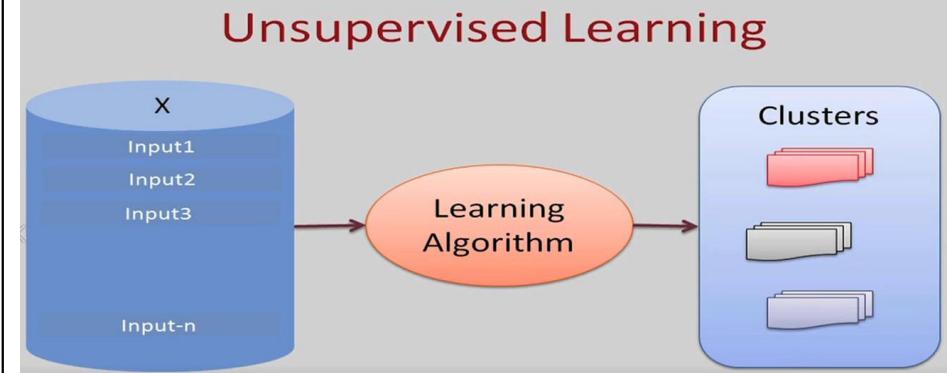


Learning Types

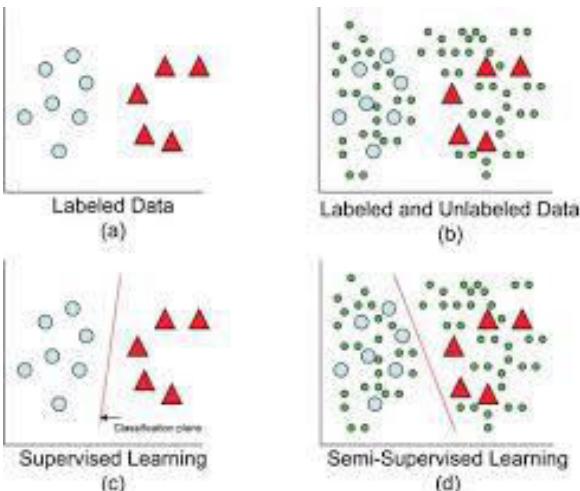


Broad types of machine learning

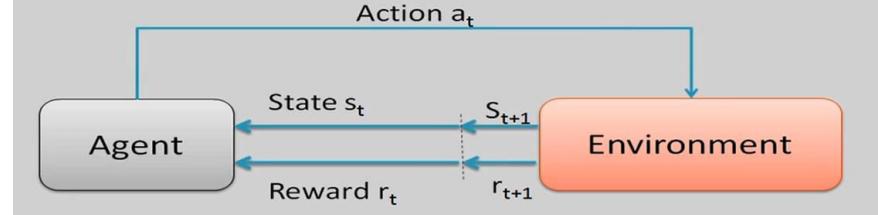
- Supervised Learning
 - X,y (pre-classified training examples)
 - Given an observation x, what is the best label for y?
- Unsupervised learning
 - X
 - Given a set of x's, cluster or summarize them
- Semi-supervised Learning
- Reinforcement Learning
 - Determine what to do based on rewards and punishments.



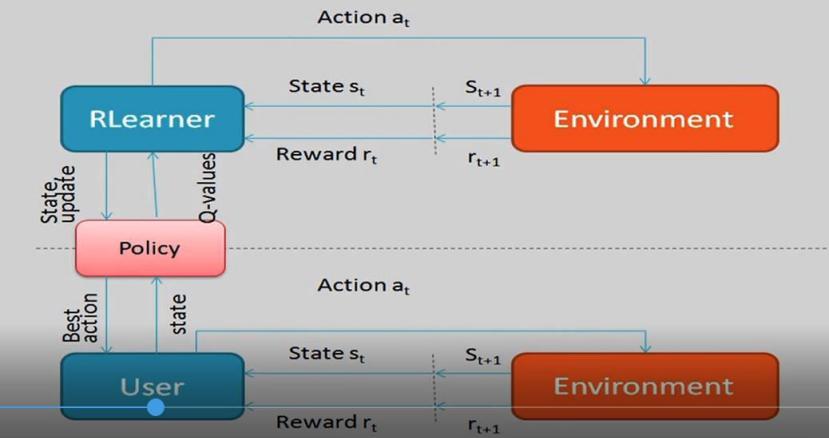
Semi-Supervised Learning



Reinforcement Learning



Reinforcement Learning



Supervised Learning

Given:

- a set of input features X_1, \dots, X_n
- A target feature Y
- a set of training examples where the values for the input features and the target features are given for each example
- a new example, where only the values for the input features are given

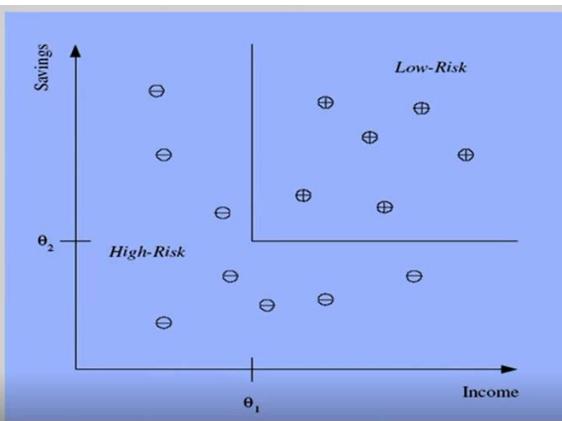
Predict the values for the target features for the new example.

- classification when Y is discrete
- regression when Y is continuous

Classification

Example: Credit scoring

Differentiating between
low-risk and high-risk
customers from their
income and savings



Discriminant: IF $income > \theta_1$ AND $savings > \theta_2$

THEN low-risk ELSE high-risk

Features

- Often, the individual observations are analyzed into a set of quantifiable properties which are called features. May be
 - categorical (e.g. "A", "B", "AB" or "O", for blood type)
 - ordinal (e.g. "large", "medium" or "small")
 - integer-valued (e.g. the number of words in a text)
 - real-valued (e.g. height)

Regression

Example: Price of a used car

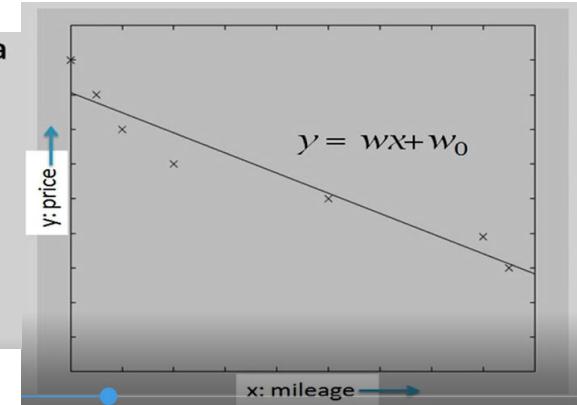
x : car attributes

y : price

$$y = g(x, \theta)$$

$g(\cdot)$ model,

θ parameters



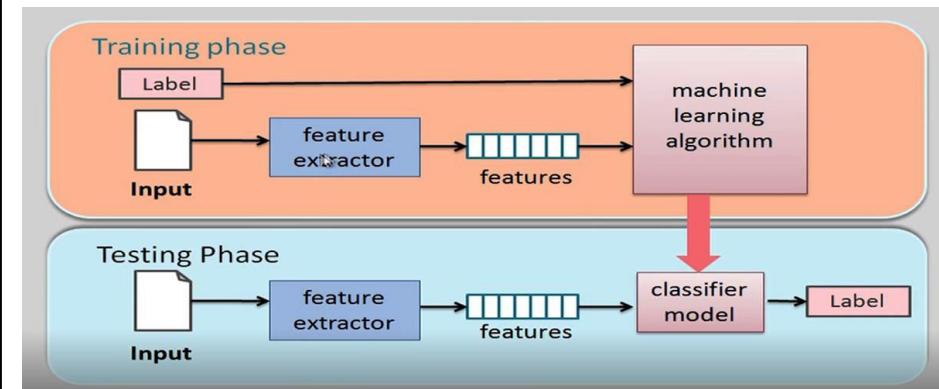
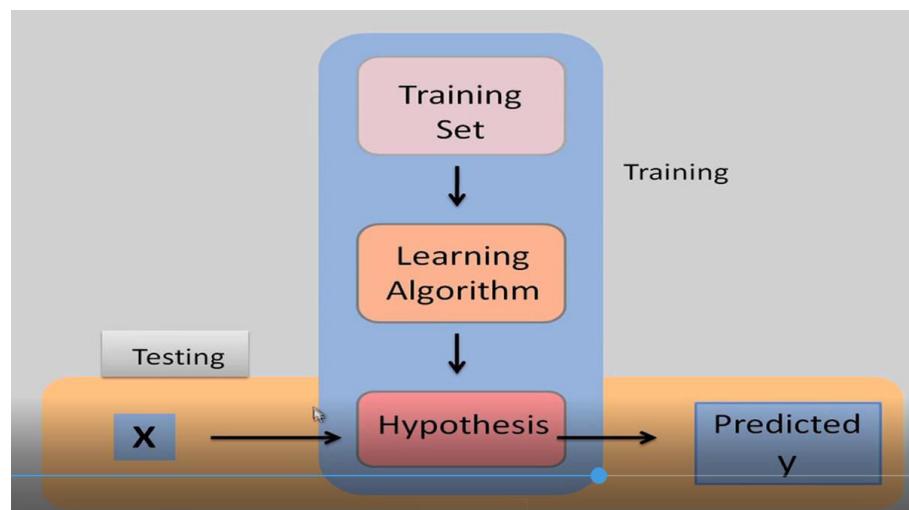
Example Data

Training Examples:

	Action	Author	Thread	Length	Where
e1	skips	known	new	long	Home
e2	reads	unknown	new	short	Work
e3	skips	unknown	old	long	Work
e4	skips	known	old	long	home
e5	reads	known	new	short	home
e6	skips	known	old	long	work

New Examples:

e7	???	known	new	short	work
e8	???	unknown	new	short	work



Classification learning

- Task T :
 - input: a set of *instances* d_1, \dots, d_n
 - an instance has a set of *features*
 - we can represent an instance as a vector $d = \langle x_1, \dots, x_n \rangle$
 - output: a set of *predictions* $\hat{y}_1, \dots, \hat{y}_n$
 - one of a fixed set of constant values:
 - $\{+1, -1\}$ or $\{\text{cancer, healthy}\}$, or $\{\text{rose, hibiscus, jasmine, ...}\}$, or ...
- Performance metric P :
- Experience E :

Classification Learning

Task	Instance	Labels
medical diagnosis	patient record: blood pressure diastolic, blood pressure systolic, age, sex (0 or 1), BMI, cholesterol	$\{-1, +1\}$ = low, high risk of heart disease
finding company names in text	a word in context: capitalized (0,1), word-after-this-equals-Inc, bigram-before-this-equals-acquired-by, ...	{first, later, outside} = first word in name, second or later word in name, not in a name
brain-human-interface	brain state: neural activity over the last 100ms of 96 neurons	{n, s, e, w, ne, se, nw, sw} = direction you intend to move the cursor

Classification learning

- Task T :
 - input: a set of *instances* d_1, \dots, d_n
 - output: a set of *predictions* $\hat{y}_1, \dots, \hat{y}_n$
- Performance metric P :
 - Prob (wrong prediction) on examples from D
- Experience E :
 - a set of *labeled examples* (x, y) where y is the true label for x
 - ideally, examples should be *sampling* from some fixed distribution D

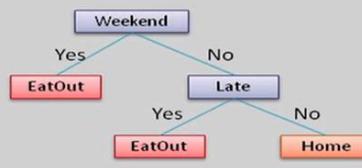
we care about performance on the distribution, not the training data

Classification Learning

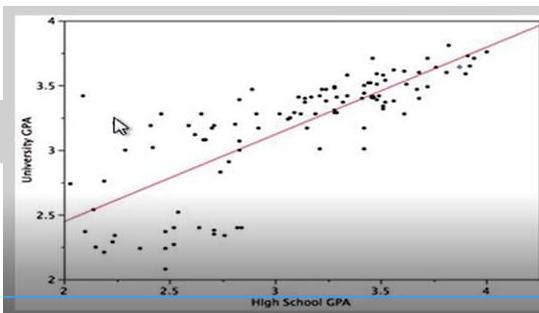
Task	Instance	Labels	Getting data
medical diagnosis	patient record: lab readings	risk of heart disease	wait and look for heart disease
finding company names in text	a word in context: capitalized, nearby words, ...	{first, later, outside}	text with company names highlighted (say by hand)
brain-human-interface	brain state: neural activity over the last 100ms of 96 neurons	{n, s, e, w, ne, se, nw, sw}	recordings of someone doing known tasks (so direction can be inferred)
image recognition	image: pixels	no house, house	hand-labeled images

Representations

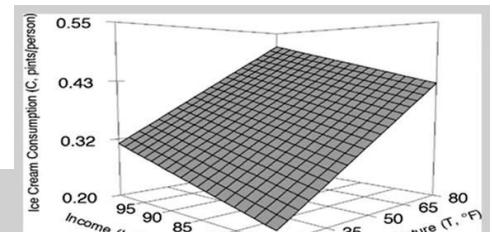
1. Decision Tree



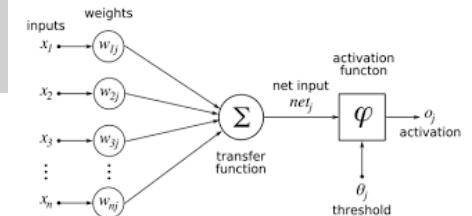
2. Linear function



3. Multivariate linear function



4. Single layer perceptron



Hypothesis Space

- One way to think about a supervised learning machine is as a device that explores a “hypothesis space”.
 - Each setting of the parameters in the machine is a different hypothesis about the function that maps input vectors to output vectors.

Terminology

- **Features:** The number of features or distinct traits that can be used to describe each item in a quantitative manner.
- **Feature vector:** n-dimensional vector of numerical features that represent some object
- **Instance Space X:** Set of all possible objects describable by features.
- **Example (x,y):** Instance x with label $y=f(x)$.

Terminology

- **Concept c:** Subset of objects from X (c is unknown).
- **Target Function f:** Maps each instance $x \in X$ to target label $y \in Y$
- **Example (x,y):** Instance x with label $y=f(x)$.
- **Training Data S:** Collection of examples observed by learning algorithm.
Used to discover potentially predictive relationships

DATA CLUSTERING

Algorithms and
Applications

Edited by
Charu C. Aggarwal
Chandan K. Reddy



CRC Press
Taylor & Francis Group
Boca Raton, London, New York
Taylor & Francis Group, an Informa business
A CHAPMAN & HALL BOOK

x	Contents
5.6 Subspace Clustering	118
5.7 Clustering Networks	120
5.8 Other Directions	123
5.9 Conclusion	124
6 Grid-Based Clustering	127
Wei Chen, Wei Wang, and Sandra Butita	
6.1 Introduction	128
6.2 The Classical Algorithms	131
6.2.1 Earliest Algorithms: GRIDCLUS and BANG	131
6.2.2 STING: The Statistical Information Grid Approach	132
6.2.3 WaveCluster: Wavelets for Grid-Based Clustering	134
6.3 Adaptive Grid-Based Algorithms	135
6.3.1 AMR: Adaptive Mesh Refinement Clustering	135
6.4 Axis-Signed Grid-Based Algorithm	136
6.4.1 NSG: Axis-Signed Grid Clustering Algorithm	136
6.4.2 ADGC: Adaptive Defect and Conger Clustering	137
6.4.3 ASGC: Axis-Shifted Grid Clustering	137
6.4.4 GDILC: Grid-Based Density-IsolLine Clustering Algorithm	138
6.5 High-Dimensional Algorithms	139
6.5.1 GLOQUE: A Class of High-Dimensional Clustering Algorithm	139
6.5.2 Variants of CLIQUE	140
6.5.3 ENCLUS: Entropy-Based Approach	140
6.5.4 3D-Grid: Grid-Based Operator Grid Partitioning	141
6.5.5 Variants of the OptGrid Approach	143
6.5.5.1 O-Cluster: A Scalable Approach	143
6.5.5.2 CBF: Cell-Based Filtering	144
6.6 Conclusions and Summary	145
7 Nonnegative Matrix Factorizations for Clustering: A Survey	149
Tao Li and Chris Ding	
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
7.3 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 Algorithmic Development	155
7.4.3 Performance 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

xi	Contents
7.5 NMF Related Enclosures	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 Other Clustering	164
7.6.6 Other Clustering Extensions	164
7.7 Conclusions	165
8 Spectral Clustering	177
Jiulu Liu and Jionwei Han	
8.1 Introduction	177
8.2 Similarity Graph	179
8.3 Unnormalized Spectral Clustering	180
8.3.1 Normalized Spectral Clustering	180
8.3.2 Unnormalized Spectral Graph Laplacian	180
8.3.3 Spectral 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 Ball-and-Cut Relaxation	186
8.5.2 Normalized Cut Relaxation	187
8.6 Random Walks View	188
8.7 Connection to Kernel K-Means and Nonnegative Matrix Factorization	189
8.8 Large Scale Spectral Clustering	191
8.9 Local Scale Spectral Clustering	192
8.10 Further Reading	194
9 Clustering High-Dimensional Data	201
Arthur Zimek	
9.1 Introduction	201
9.2 The "Curse of Dimensionality"	202
9.2.1 Different Aspects of the "Curse"	202
9.2.2 Countermeasures	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
9.3.2 Search Space	207
9.3.2.1 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

xii	Contents
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
Charu C. Aggarwal	
10.1 Introduction	231
10.2 Methods Based on Partitioning Representatives	233
10.2.1 TREC-STREAM Algorithm	233
10.2.2 SubSpace: 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 ChiStream	237
10.3 Density-Based Streaming Clustering	239
10.3.1 DenStream: Density-Based Microclustering	240
10.3.2 Grid-Based Streaming Algorithm	241
10.3.2.1 D-Stream Algorithm	241
10.3.2.2 Other Grid-Based Algorithms	242
10.4 Probabilistic Stream Clustering	243
10.5 Clustering High-Dimensional Streams	243
10.5.1 The HPSSTREAM Method	244
10.5.2 Other High-Dimensional Streaming Algorithms	244
10.6 Clustering Uncertain Data Streams	245
10.6.1 Clustering Binary Data Streams with k-Means	245
10.6.2 The StreamChiCD Algorithm	245
10.6.3 Massive Domain Clustering	246
10.7 Other Techniques for Stream Clustering	249
10.8 Other Searches 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
Hanghang Tong and Yung Kong	
11.1 One-Pass Clustering	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 CUR: Fighting with the Irregular Clusters	263
11.3 Parallel and Distributed 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 Grid-Based Framework	268
11.4.2 BRICS: Density-Based Clustering	269
11.4.3 ParMeTIS: Graph Partitioning	269
11.4.4 PKMeans: K-Means with MapReduce	270
11.4.5 DiCoC: Co-Clustering with MapReduce	271
11.4.6 BoW: Subspace Clustering with MapReduce	272
11.5 Conclusion	274

xiii	Contents
12 Clustering Categorical Data	277
Bill Andriopoulos	
12.1 Introduction	278
12.2 Goals of Categorical Clustering	279
12.2.1 Clustering Mixed Data	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 Jaccard/Bonferroni Similarity Measures	284
12.4 Descriptions of Algorithms	284
12.4.1 Partition-Based Clustering	284
12.4.1.1 F-Modes	284
12.4.1.2 L-Projections (Mixed Categorical and Numerical)	285
12.4.1.3 L-Subspaces	286
12.4.1.4 Squeezes	286
12.4.1.5 COOLCAT	286
12.4.2 Hierarchical Clustering	287
12.4.2.1 Agglomerative Clustering	287
12.4.2.2 CORWEB	288
12.4.2.3 LIMBO	289
12.4.3 Density-Based Clustering	290
12.4.3.1 Projected (Subspace) Clustering	290
12.4.3.2 STIRR	290
12.4.3.3 CLICKS	291
12.4.3.4 STIRR	291
12.4.3.5 CLOPE	292
12.4.3.6 HDBSCAN: Hierarchical Density-Based Clustering	292
12.4.3.7 MULTIC: 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 L-Subspaces (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
David C. Anastasiou, Andrea Tagarelli, and George Karvounis	
13.1 Introduction	306
13.2 Modeling a Document	306
13.2.1 Preliminaries	306
13.2.2 Vector Space Model	307
13.2.3 Alternate Document Models	309
13.2.4 Dimensionality Reduction for Text	309
13.2.5 Clustering Extensions	310
13.3 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

Contents	xviii
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 Conclusion	530
21 Alternative Clustering Analysis: A Review	\$38
<i>James Bailey</i>	
21.1 Introduction	535
21.2 Technical Preliminaries	537
21.3 Multiway 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 Naïve	539
21.3.2.2 Meta Clustering	539
21.3.2.3 Eigenvectors of the Laplacian Matrix	540
21.3.2.4 Decoupled k-Means and Convolutional EM	540
21.3.2.5 CAMI	540
21.3.3 Guided Generation with Constraints	541
21.3.3.1 NACI	541
21.3.3.2 Constraint Optimization Approach	542
21.3.4 Orthogonal Transformation Approaches	543
21.3.4.1 Orthogonal Views	543
21.3.4.2 AIDE	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 CIB	545
21.3.6 Connections to Multiview Clustering and Subspace Clustering	545
21.3.7 Future Research Issues	547
21.3.8 Summary	547
22 Cluster Ensembles: Theory and Applications	\$51
<i>Jerudeep Ghosh and Asyan Acharya</i>	
22.1 Introduction	551
22.2 The Cluster Ensemble Problem	554
22.3 The 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 Ensemble (BCE)	558
22.4.1.3 Nonparametric Bayesian Cluster Ensemble (NPBCE)	559
22.4.2 Pairwise Similarity-Based Approaches	560
22.4.2.1 Methods Based on Ensemble Co-Association Matrix	560

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

Contents	xix
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 Application of Consensus Clustering	564
22.5.1 Geospatial Data Analysis	564
22.5.2 Image Segmentation	564
22.6 Concluding Remarks	566
23 Clustering Validation Measures	\$71
<i>Hui Xiong and Zhongzhou 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 Internal Clustering Validation Measures	575
23.2.2.1 K-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 The Impact of Defective Measures	577
23.2.3 Measure Normalization	577
23.2.3.1 Normalizing the Measures	578
23.2.3.2 The DCV Criterion	581
23.2.3.3 The Impact 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 The Impact of Distance	589
23.2.5 Internal Clustering Validation Measures	589
23.2.5.1 An Overview of Internal Clustering Validation Measures	589
23.2.5.2 Understanding of Internal Clustering Validation Measures	592
23.2.5.2.1 The Impact of Monotonicity	592
23.2.5.2.2 The Impact of Noise	593
23.2.5.2.3 The Impact of Density	594
23.2.5.2.4 The Impact of Subclusters	595
23.2.5.2.5 The Impact of Skewed Distributions	596
23.2.5.3 The Impact of Arbitrary Shapes	598
23.2.5.4 Properties of Measures	600
23.2.6 Summary	601
24 Educational and Software Resources for Data Clustering	\$67
<i>Churni Garg 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
24.2.3 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

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, image, bioinformatics, and web mining. Clustering has been applied in a wide variety of 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 the main aspects of clustering. The first set of chapters will focus on core methods for data clustering. These include hierarchical clustering, density-based 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 focus on different applications of clustering problems. This includes the sensitivity 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 theoretical, empirical, and with the help of various aspects of machine learning. The book will focus on the core methods for data clustering. These include hierarchical clustering, density-based 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 streaming or uncertain data. Clustering has been addressed by a number of different approaches such as pattern recognition, data mining and machine learning. In some cases, the work to fit the different communities tends to be disjointed 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 complex and specialized data clustering methods. Special attention will also be paid to more recent problem domains such as graph and social networks.

- **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, and uncertain data. These chapters also discuss semi-supervised 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. Orlitz. He has since worked in the field of performance analysis, databases, and data mining. He has published over 200 papers in international conferences, journals, and books, and has 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 three been designated a Master Inventor at IBM. He is a recipient of an IBM Corporate Award for his technical work in 2002, a recipient of a research award of the Research Division of 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 ICDE Conference (2007), the WWW Conference (2008), and the ACM SIGKDD 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 active 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, "contributions to knowledge discovery and data mining techniques", and a life-member of the ACM.

S. K. Dasgupta is an Associate 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 and 100+ conference proceedings. 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 Agothic
Relancy, LLC
Saint Louis Park, Minnesota

Mohammad Al Hasan
Indiana University - Purdue University
Indianapolis, Indiana

Salem Alejani
Arizona State University
Tempe, Arizona

David C. Anastasiu
University of Minnesota
Minneapolis, Minnesota

Bill Andreadopoulos
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-cha-cha 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	Jiliang Tang Arizona State University Tempe, Arizona
Dimitris Karatzos University of Athens Athens, Greece	Hanqiu Tang IBM T.J. Watson Research Center Yorktown Heights, New York
Tao Li Florida International University Miami, Florida	Gece Trajcevski Northwestern University Evanston, Illinois
Zhongming Li Rutgers University New Brunswick, New Jersey	Min-Hsuan Tsai University of Illinois at Urbana-Champaign Urbana, Illinois
Huan Liu Arizona State University Tempe, Arizona	Shen-Fu Tsai Microsoft Inc. Redmond, Washington
Jiali Lin University of Illinois at Urbana-Champaign Urbana, Illinois	Bhanukiran Vimalsuri Wayne State University Detroit, Michigan
Srinivasan Parthasarathy The Ohio State University Columbus, Ohio	Wei Wang University of California Los Angeles, California
Guo-Jun Qi University of Illinois at Urbana-Champaign Urbana, Illinois	Hui Xiong Rutgers University New Brunswick, New Jersey
Chandan K. Reddy Wayne State University Detroit, Michigan	Mohammed J. Zaki Rensselaer Polytechnic Institute Troy, New York
Andrea Tagarelli University of Calabria Arcavacata di Rende, Italy	Arthur Zimek University of Alberta Edmonton, Canada

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 on the specific model used. For example, a generative model may define similarity on the basis of a probabilistic generative model, 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 processes 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, 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 wide 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. An example 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 sensor data, and trajectory data. Key trends and events in the data can be discovered with the use of clustering methods.
- **Multimodal Data Analysis:** There are various types of data 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

Chapter 1

An Introduction to Cluster Analysis

Charu C. Aggarwal
IBM T.J. Watson Research Center
Yonkers 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-Based 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	Kernel-Based 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	Storage 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 Numerical 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 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

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 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 they are well suited in different data 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-Centred:** 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 type will dictate the choice of methods suitable 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 quality of the clusters.

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 algorithm. It is often the case that the features in the clusters may sometimes be more than other. Therefore, it is often helpful to utilize a pre-processing 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 reduce the feature selection effect. The number of the former is called *reducibility*, while the number 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 data examples prior knowledge with respect to a given set of features ... can be used ... in conjunction with a feature selection strategy to define dimensions that are likely to be feasible to prune off many dimensions without at the same time incurring a substantial loss of information" [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 set is used to estimate the parameters in such a way that they have a *maximum likelihood fit* to the generative model. Generative models are based on the joint probability distributions (or 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 distribution $\mathcal{G}_1 \dots \mathcal{G}_k$ 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 \mathcal{G}_r .

The probability distribution \mathcal{G}_r is picked from a set of different possibilities. Note that this generates a probability distribution over a set of several distributions, as the prior probabilities and the model parameters for each distribution \mathcal{G}_i . 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 \mathcal{G}_i . Some examples are as follows:

- For numerical data, a Gaussian mixture model may be used for \mathcal{G}_i in order to model each component \mathcal{G}_i . Such a model is discussed in detail in Chapter 3.
- For categorical data, a Bernoulli model may be used for \mathcal{G}_i in order to model the generation of the discrete values.
- For sequence data, a Hidden Markov Model (HMM) may be used for \mathcal{G}_i 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 the prior probabilities for each cluster are equal 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 underlying data space. Examples include the estimation of the mean and the covariance probabilities in terms of the euclidean 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 representatives are updated. This is a very similar process to the k -means algorithm, except that we 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 euclidean distance is used in order to create the clusters. The k -Means approach is one of the simplest and most classical methods for data clustering [66] 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 much more stable to noise outliers than the k -Means approach [67], but it is usually less sensible to extremes 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 overlapping and confusion in the literature, it should be clear that the k -Means and k -Medians methods should be considered as two distinct techniques. Throughout the several chapters of this book, the k -Medians approach discussed is really a k -Medoids approach, though we have chosen to be consistent with the specific research paper which is being described. Nevertheless, it is important to note the overlapping of this term in order to avoid confusion.

– *k -Medoids*: In these methods, partitioning representatives are drawn 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 move the quality of these representatives. In this 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 -Medians 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 then merge them in order to create larger and larger clusters. 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 of points in two clusters is used, a measure of 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 the number of branches of the tree and the depth of the tree 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 balance of the node depths and the number of data points in the branches. For example, in a tree with 1000 nodes, 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 equally 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 text documents [83] or biological sequences [84]. In some cases, the distance function and data quality, the quality of the distance functions reduces because of many irrelevant dimensions [41] 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 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 smooth function that represents the density of the data space. In a region of high density of granularities, the data points are said to be "together" or "dense". Grid-based methods are 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 representing the data 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 they explore the data space at a high level of granularity, which can be recommended for certain kinds of data. There are two main classical methods for density-based methods and grid-based methods: 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 nonEuclidean 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 feature selection and clustering, in that they attempt to use the clustering process to better dimensionality 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 proximities analysis, as opposed to the rows. This is where the question arises, as to whether it is possible to perform two steps simultaneously, by clustering rows and columns 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 matrix factorization are also called dimensionality reduction based on the same general principle. 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 generative model for document clustering can be considered a model of assignment of documents to topics, 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 using the hard choices of point and dimension selection by nonparametric methods. Nevertheless, many special cases of generative data types have been learned successfully with generative models.

A particularly common dimensionality reduction method is that of *Probabilistic Latent Semantic Indexing* (PLSI). This method is also sometimes referred to as *Probabilistic Latent Semantic Indexing (PLSI)*. In topic modeling clusters are 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 sparse matrix of the probabilities of word-document assignments in a topic-specific manner 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-

on the rows (documents) and columns (words). This general principle, when applied to sparse nonnegative 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 a 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 a dual task of mapping these objects into a coordinate space, while preserving the dimensionality. Therefore, spectral methods are also useful for performing clustering on arbitrary objects such as nodes 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, which is an expensive process, takes a very small number of these operations in practice. 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 X_i and X_j , the *heat kernel* defines a similarity matrix W of weights:

$$w_{ij} = \exp(-\|\overrightarrow{X_i} - \overrightarrow{X_j}\|^2/t) \quad (1.3)$$

Here t is a user-defined parameter. Furthermore, if w_{ij} value of W_{ij} is set to 0 if the distance between $\overrightarrow{X_i}$ and $\overrightarrow{X_j}$ 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 solving eigenvalues of this graph, which consists of pairs which are weakly connected by edges representing similarity. Hence, spectral methods be considered as a graph-based technique for clustering of any kind 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 examples include the minimum k -nearest neighbor graph, simply the binary graph in which the distances are less than a user-specified threshold t .

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 eigenvectors 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 euclidean distances on this line. Therefore, it is undesirable for data points which are very similar to be mapped

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

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

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 $\mathbf{y} = (y_1 \dots y_n)$. The objective function O can be rewritten as follows:

$$\mathbf{y}^T \mathbf{y} - \mathbf{y}^T L \mathbf{y} \quad (1.5)$$

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:

$$\mathbf{y}^T \mathbf{y} = 1 \quad (1.6)$$

Note that the use of the matrix D provides different weights to the data items, because it is assumed that the data items are mapped to different dimensions are mapped in the same relative order. This optimization formulation is a generalized eigenvalue problem, in which the value of $\mathbf{y}^T \mathbf{y}$ is optimized by selecting the smallest eigenvalue for which the generalized eigenvector relationship $L \mathbf{y} = \lambda D \mathbf{y}$ is satisfied. In practice however, the smallest generalized eigenvector corresponds to the trivial solution, where \mathbf{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 corresponds to the most informative eigenvector in a more meaningful way.

This model can be generalized to determine 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 eigenvectors $\overrightarrow{\mathbf{y}_1}, \dots, \overrightarrow{\mathbf{y}_n}$ (of which the first is trivial), with corresponding eigenvalues $\lambda_1 > \lambda_2 > \dots > \lambda_n$. Let the corresponding vector representation of the data points along each eigenvector be denoted by $\overrightarrow{\mathbf{y}}_i$. What do the small and large magnitude eigenvectors intuitively represent in the transformed space? By using the ordering of the data along a small magnitude eigenvector to create a cut, the weight of the edges across the cut is likely to be small. Thus, that represents a cluster in the space of data. At the same time, a large magnitude eigenvector will provide a lower dimensional (possibly nonlinear) shape of the data, though this comes at a rather high cost [78]. A 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 number of dimensions of the data attributes of different parts of the data. This leads to numerous challenges in many of the main problems such as clustering, nearest neighbors search, and outlier analysis. It should be noted that many of the algorithms mentioned so far are not problem 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 more features are added and correlated, 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 known as high-dimensional noise and concentration. 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. 508. These noise and concentration effects are therefore a problematic symptom of (locally) irrelevant, uninformative or noisy attributes, which have additive effects that statistically significantly full-dimensional algorithms are led to ignore in the notion of distances [44]. In particular, it would seem odd that data mining algorithms should behave poorly with increasing dimensionality, if learning a query point’s neighbors with a large 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 of locally irrelevant attributes. Carefully designed distance functions can leverage the greater inter-point distances to improve the 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 is performed specifically to different localities of the data. Some projected methods even refer to themselves 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 $T_1 \dots T_k$ such that the points in C_i cluster well in the subspace represented by T_i . Thus, these methods are a form of local clustering, where they are used to find the most relevant clusters in the underlying data. Note that the subspace T_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 projected clustering has increased since 2000, survey on this topic can be found in [65, 54]. An overview of algorithms for high-dimensional data can be provided in Chapter 9.

It should also be pointed out that while pairwise distances are often 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 redundancy in representing different pairwise instances. Therefore, even when there is significant noise and concentration effects, one will often be able to use the distance matrix as a standard Euclidean 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

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 requirement of real-time analysis, but robustness and energy efficiency 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 frame to capture the evolving aspects of the underlying data stream. Streaming data can also be used for other data types such as discrete data, massive 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 some 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 amount of data can make a distributed system hard to be used in order to store all the data, and distributed processing becomes increasingly difficult in order to maintain 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 systems, it is desirable to exchange information between nodes in an efficient and cost-effective way. For an application programmer, this sometimes creates 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 model for analysis of large amounts of data, especially when the size of the data is too large to be linearly computable statistician functions over the entire 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 input data, scheduling the program’s execution across a set of machines, handling machine failures, and managing the movement of data between machines. The application programmer need only express how to parallelize and distribute his program*” 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* can be found in [1]. The key idea is to use a *Map* function to distribute data among the different machines, and then use an *Reduce* function to sum up much smaller data (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

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 too 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 the data points in the space. These findings 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 expected to work well in high-dimensional space, off-the-shelf methods can be used to mitigate concentration effects when they do retain sufficient information collectively when used carefully in conjunction with dimensioning 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 cluster 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 this case, the data is collected sequentially over time, and the overall system, 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 requirements for processing the objects in the database. Therefore, algorithms which use sequential access of the data, namely those that sequentially access the same records, often prove 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 the case are flat partitioning algorithms, which are essentially scans of the data 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-means approach. Note that the k-means approach can still be computationally quite intensive because each iteration requires trying out new partitioning representatives through an exchange of data points. If a data set is very large, then it may take many iterations to converge. A smaller number of passes over the data set is needed if 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

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 high-level interface which is independent of very specific machine 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 systems 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 categories are binary.

Categorical data sets lead to several challenges for clustering algorithms:

- When the algorithms depend upon the use of a similarity or distance function, the standard L_2 metric can no longer be used. New similarity measure 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 data is mixed, then the problem becomes more difficult because the different attributes now need to be treated in a heterogeneous way, and the similarity functions need to explicitly account for this heterogeneity.

It should be noted that generative 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, it is possible to cluster categorical data. Specifically, hierarchical clustering is one class of models which can be used with virtually any type of data. As an example, the L_1 distance can be defined. The downside is that the methods scale with the square of the similarity matrix size. Generative model 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 stored in a bag-of-words format, in which the specific ordering of abstract words 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, 70] use distance based metrics. Since the data is high dimensional and nonnegative, one needs to take care 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 as versions of models such as *LSI* [31] and *SVD* [53]. In addition, it is possible to 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 for text clustering are known to closely relate to spectral clustering [88]. In recent years, the popularity of social media has also led 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 *Facebook* 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 textual 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 multimodal, while audio and video data are unimodal. The context in which the data is collected 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

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 can be found in [58]. One key problem which arises in the context of discrete sequences is that the coordinate 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 coordinate 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 CLUSSEQ [10].

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 hidden Markov models (HMMs) for sequence clustering, as discussed in Chapter 7 [69]. 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 can be represented as a similar 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 for the objects, graph clustering can be used in order to perform the analysis. Graph clustering has been applied extensively in the social network context, 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 have also involved 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 these methods often lead to high-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 partition. Other objective functions for creating cuts of the sort, 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.

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 a form of data which is often referred to as sensor data, stock market data, or 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 significance difference in problem definition 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 used or shape-based offline analysis.

- 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 over time. In this case, the clusters formed by different stocks need to be updated continuously, based on their mutual regression coefficients. Some examples of these methods include the MUSCLES approach [85] 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 given ECG time series, collected from a patient, the specific time of when the series was collected is unimportant, 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 dimension. As in the previous case, the design of the distance function [42] holds the key to the effective use of the approach.

A particularly 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 data can be viewed as a collection of multiple time series, all of which are related to the same series. It is possible to perform these analyses 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 focused 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.

• **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 corresponding solvers for these problems can be used to create a multi-level graph clustering 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 adjacencies or to the edge weights. This can be used to gain certain different kinds of insights. It is also relatively easy to augment the matrix with structure 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 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-value uncertainty or in the form of a *possible world* scenario, in which only particular subsets of attributes can be true in any one of a group of worlds [3]. The key idea here is that 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 useful for the algorithm. This is particularly true in the context of streaming data in 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 suits the needs of the application? There are numerous ways in which these issues can be explored. The explanations could be through interactive visualization and human interaction, external knowledge-based supervision,

especially 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 use 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 known systems for exploring high-dimensional data is the *HiD-Eye* method [144], which explores the differences of the data in order to determine clusters in different feature-specific views of the data. Another well-known technique is the *PCCLUS* 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 user can then choose a specific view and count the number of clusters, 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 a diagnostic step in many clustering algorithms, it becomes important to have a good list of clusters that may suit that application. The sufficiency 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 subjective analysis also implies that it is difficult to refer to one particular clustering as significantly better than another. It is clear 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 clustering solution. He may want to use a clustering algorithm such as *k-means* to explore 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, 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 density, but a density-based method may be superior to a *k-means* algorithm in terms of the estimated density of the clusters. This is because there is a circularity in using particular criteria 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*) due to 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 evaluated by comparing 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, feature-based methods, and constraint-based methods. The problems of clustering have close relationships with 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 integral part of feature selection/dimensionality reduction methods with the use of large amounts of data through an ever-increasing number of ways. The use of distributed 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

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 elements of the appropriate categories are provided as (some of the) seeds to a regular 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 setting, as it relates to the underlying clusters. The reader is referred to [18] for a detailed discussion of semisupervised clustering. An excellent discussion on constrained clustering algorithms may be found in [18], 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 clustering are possible. When supervision is available, the better clustering of properties in the clustering process can be used to perform better analysis, especially from the perspective of interpretability. These are referred to as *informative 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 is often referred to as *multiview clustering*.

One of the most used method for multiview clustering is to run clustering algorithms multiple times and then examine the different clustering solutions to determine 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 that are orthogonal to each other. Thus, by applying a 2-means algorithm on 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. Eigenvectors are orthogonal eigenvectors, where the clusters 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, ensemble-based 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 all clustering algorithms.

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 often is 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 density, but a density-based method may be superior to a *k-means* algorithm in terms of the estimated density of the clusters. This is because there is a circularity in using particular criteria 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*) due to 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 evaluated by comparing 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.

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 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. Hinneberg, 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 spaces. *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.

- 24 *Data Clustering: Algorithms and Applications*
- [18] S. Basu, J. 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? *KDD Conference*, 2001.
- [22] D. Blei, A. Ng, and M. Jordan. Latent Dirichlet allocation. *Journal of Machine Learning Research*, 3:993–1022, 2003.
- [23] T. 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. Traas Jr., A. Traas, 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://taana.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. Camasta, 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.

- 26 *Data Clustering: Algorithms and Applications*
- [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] J. Liu. Clustering of time series data—A survey. *Pattern Recognition*, 38(11):1857–1874, 2005.
- [61] J. Liu and H. Motoda. *Computational Methods of Feature Selection*. Chapman and Hall (CRC Press), 2007.
- [62] J. Liu, R. Jin, C. Aggarwal, and Y. Shen. Reliable clustering on uncertain graphs. *ICDM Conference*, 2012.
- [63] U. von Luxburg. A tutorial on spectral clustering. *Statistics and Computing*, Springer, 2007.
- [64] A. L. 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] J. Mourtaghi. A survey of recent advances in hierarchical clustering algorithms. *The Computer Journal*, 26(4):354–359, 1983.
- [66] B. 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. Rabine. 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, 2008.
- [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] J. W. Silverman. *Density Estimation for Statistics and Data Analysis*. Chapman and Hall, 1986.
- [75] H. Schütze 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.

- 25 *An Introduction to Cluster Analysis*
- [77] G. Gan, C. Ma, and J. Wu. *Data Clustering: Theory, Algorithms, and Applications*. SIAM, Society for Industrial and Applied Mathematics, 2007.
- [78] V. Ganti, J. Gehrke, and R. Ramakrishnan. CACTUS-clustering categorical data using summaries. *ACM KDD Conference*, 1999.
- [79] D. Gibson, J. Kleberg, and P. Raghavan. Clustering categorical data: An approach based on Dynamical Systems. *VLDB Conference*, 1998.
- [80] S. Guha, R. Rastogi, and K. Shim. ROCK: A robust clustering algorithm for categorical attributes. *Information Systems*, 25(5):345–366, 2000.
- [81] S. Guha, R. Rastogi, and K. Shim. CURIE: An efficient clustering algorithm for large databases. In *ACM SIGMOD Conference*, 1998.
- [82] D. Gunopulos and G. Das. Time-series similarity measures, and time series indexing. *ACM SIGMOD Conference*, 2001.
- [83] A. Hinneburg. What is the nearest neighbor in high dimensional spaces. *VLDB Conference*, 2000.
- [84] A. Hinneburg, D. Keim, and M. Waagenknecht. Hd-tree: Visual mining of high-dimensional data. *IEEE Computer Graphics and Applications*, 19(2):31–31, 1999.
- [85] T. Hofmann. Probabilistic latent semantic indexing. *ACM SIGIR Conference*, 1999.
- [86] A. Jain. Data clustering: 50 years beyond k-means. *Pattern Recognition Letters*, 31(8):651–666, 2010.
- [87] A. Jain and R. Dubes. *Algorithms for Clustering Data*. Prentice Hall, 1988.
- [88] A. Jain, M. Murty, and P. Flynn. Data clustering: A review. *ACM Computing Surveys (CSUR)*, 31(3):264–323, 1999.
- [89] 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.
- [90] I. Jolliffe. *Principal Component Analysis*. Springer, 2002.
- [91] D. R. Karger. Random sampling in cut, flow, and network design problems. *STOC*, pp. 648–657, 1994.
- [92] L. Kaufman and P. Rousseeuw. *Finding Groups in Data: An Introduction to Cluster Analysis*. Wiley-Interscience, 2005.
- [93] G. Karypis and V. Kumar. Multilevel k-way partitioning scheme for irregular graphs. *Journal of Parallel and Distributed Computing*, 48(1):96–129, 1998.
- [94] 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.
- [95] D. D. Lee and H. S. Seung. Learning the parts of objects by non-negative matrix factorization, *Nature*, 401:788–791, 1999.
- [96] B. W. Kernighan and S. Lin. An efficient heuristic procedure for partitioning graphs. *Bell System Technical Journal*, 49:291–307, Feb 1970.

- 27 *An Introduction to Cluster Analysis*
- [97] 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.
- [98] J. Tenenbaum, B. de Silva, and J. Langford. A global geometric framework for nonlinear dimensionality reduction. *Science*, 290(5500):2319–2323, 2000.
- [99] P. Willett. Recent trends in hierarchical document clustering: A critical review. *Information Processing & Management*, 24(5):577–597, 1988.
- [100] R. Xu and D. Wunsch. Survey of clustering algorithms. *IEEE Transactions on Neural Networks*, 16(3):645–678, 2005.
- [101] R. Xu and D. Wunsch. *Clustering*. Wiley-IEEE Press, 2008.
- [102] F. Wang and J. Sun. Distance metric learning in data mining. *SDM Conference (Tutorial)*, 2012.
- [103] W. Wang, J. Yang, and R. Muntz. Sting: A statistical information grid approach to spatial data mining. *VLDB Conference*, 1997.
- [104] T. White. *Hadoop: The Definitive Guide*. Yahoo! Press, 2011.
- [105] J. Yang and W. Wang. CLUSEQ: Efficient and effective sequence clustering. *ICDE Conference*, 2003.
- [106] 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.
- [107] T. Zhang, R. Ramakrishnan, and M. Livny. BIRCH: An efficient data clustering method for very large databases. In *ACM SIGMOD Conference*, 1996.
- [108] Y. Zhao, G. Karypis, and U. Fayyad. Hierarchical clustering algorithms for document datasets. *Data Mining and Knowledge Discovery*, 10(2):141–168, 2005.
- [109] Y. Zhou, H. Cheng, and J. Yu. Graph clustering based on structural/attribute similarities. *Proc. VLDB Endow.*, 2(1):718–729, 2009.
- [110] Y. Zhu and D. Shah. StatStream: Statistical monitoring of thousands of data streams in real time. *VLDB Conference*, pages 358–369, 2002.

Concept Learning

General-to-Specific Learning Most

Most General Hypothesis: $h = <?, ?, ?, ?, ?, ?, ?>$

Every day Tom has enjoy i.e., Only positive examples.

Most Specific Hypothesis: $h = <\emptyset, \emptyset, \emptyset, \emptyset, \emptyset, \emptyset, \emptyset>$

Concept Learning as Search:

- Concept learning can be viewed as the task of searching through a large space of hypothesis implicitly defined by the hypothesis representation.
- The goal of the concept learning search is to find the hypothesis that best fits the training examples.

General-to-Specific Learning:

Example	Sky	AirTemp	Humidity	Wind	Water	Forecast	EnjoySport
1	Sunny	Warm	Normal	Strong	Warm	Same	Yes
2	Sunny	Warm	High	Strong	Warm	Same	Yes
3	Rainy	Cold	High	Strong	Warm	Change	No
4	Sunny	Warm	High	Strong	Cool	Change	Yes

$$h_1=(\text{Sunny}, ?, ?, \text{Strong}, ?, ?, ?)$$
$$h_2=(\text{Sunny}, ?, ?, ?, ?, ?, ?)$$

*h2 is more general than h1.
*h2 imposes fewer constraints on the instance than h1.

Version Space:

The set of all valid hypotheses provided by an algorithm is called version space (VS) with respect to the hypothesis space H and the given example set D.

Candidate-Elimination Algorithm:

- * The Candidate-Elimination algorithm finds all describable hypotheses that are consistent with the observed training examples.
- * Hypothesis is derived from examples regardless of whether x is positive or negative example

Initialize G to the set of maximally general hypotheses in H

Initialize S to the set of maximally specific hypotheses in H

For each training example d , do

- If d is a positive example
 - Remove from G any hypothesis inconsistent with d
 - For each hypothesis s in S that is not consistent with d
 - Remove s from S
 - Add to S all minimal generalizations h of s such that
 - h is consistent with d , and some member of G is more general than h
 - Remove from S any hypothesis that is more general than another hypothesis in S
- If d is a negative example
 - Remove from S any hypothesis inconsistent with d
 - For each hypothesis g in G that is not consistent with d
 - Remove g from G
 - Add to G all minimal specializations h of g such that
 - h is consistent with d , and some member of S is more specific than h
 - Remove from G any hypothesis that is less general than another hypothesis in G

Candidate-Elimination Algorithm:

- It employs a much more compact representation of the version space.
- In this the version space is represented by its most general and least general members (Specific).
- These members form general and specific boundary sets that delimit the version space within the partially ordered hypothesis space.

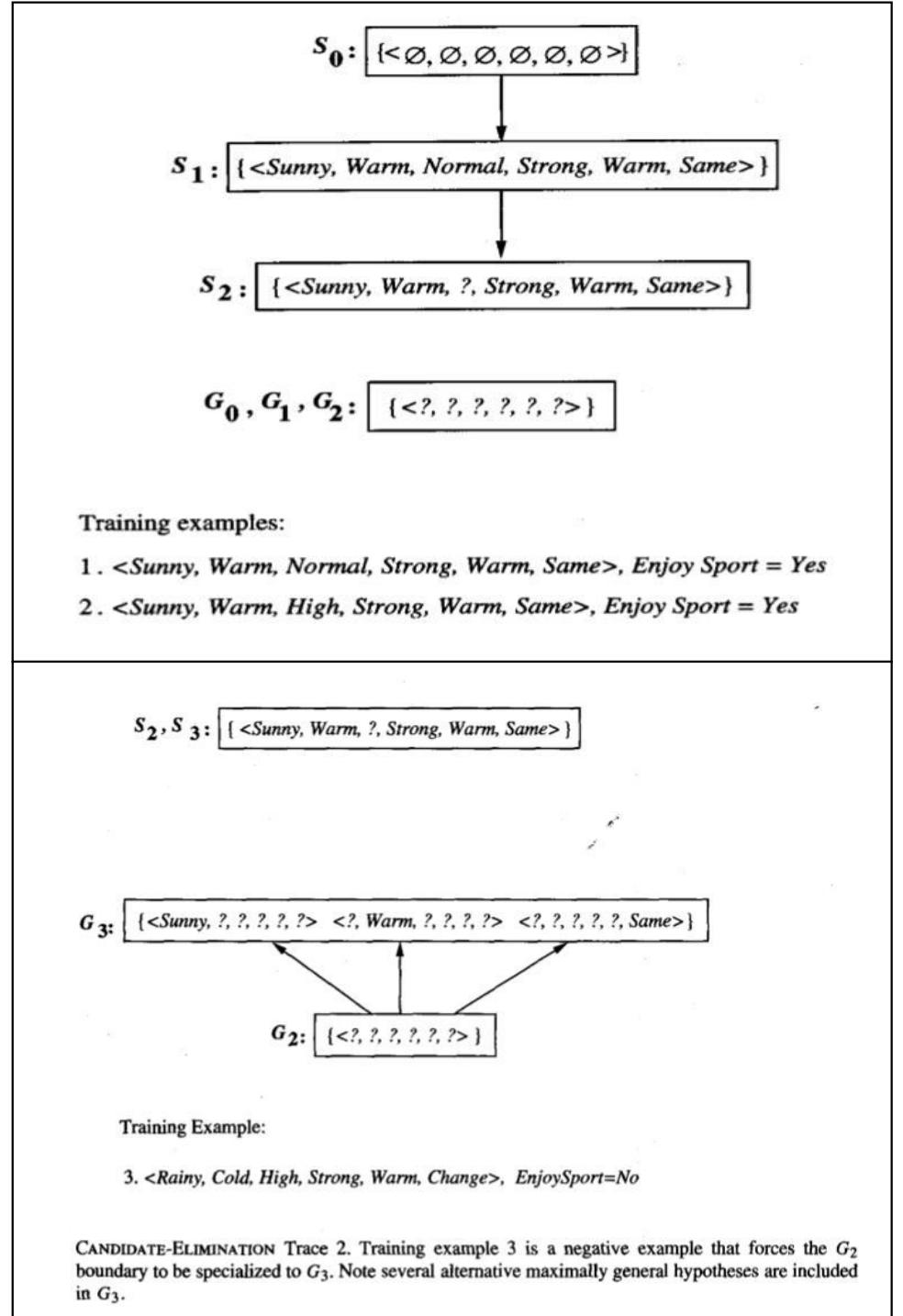
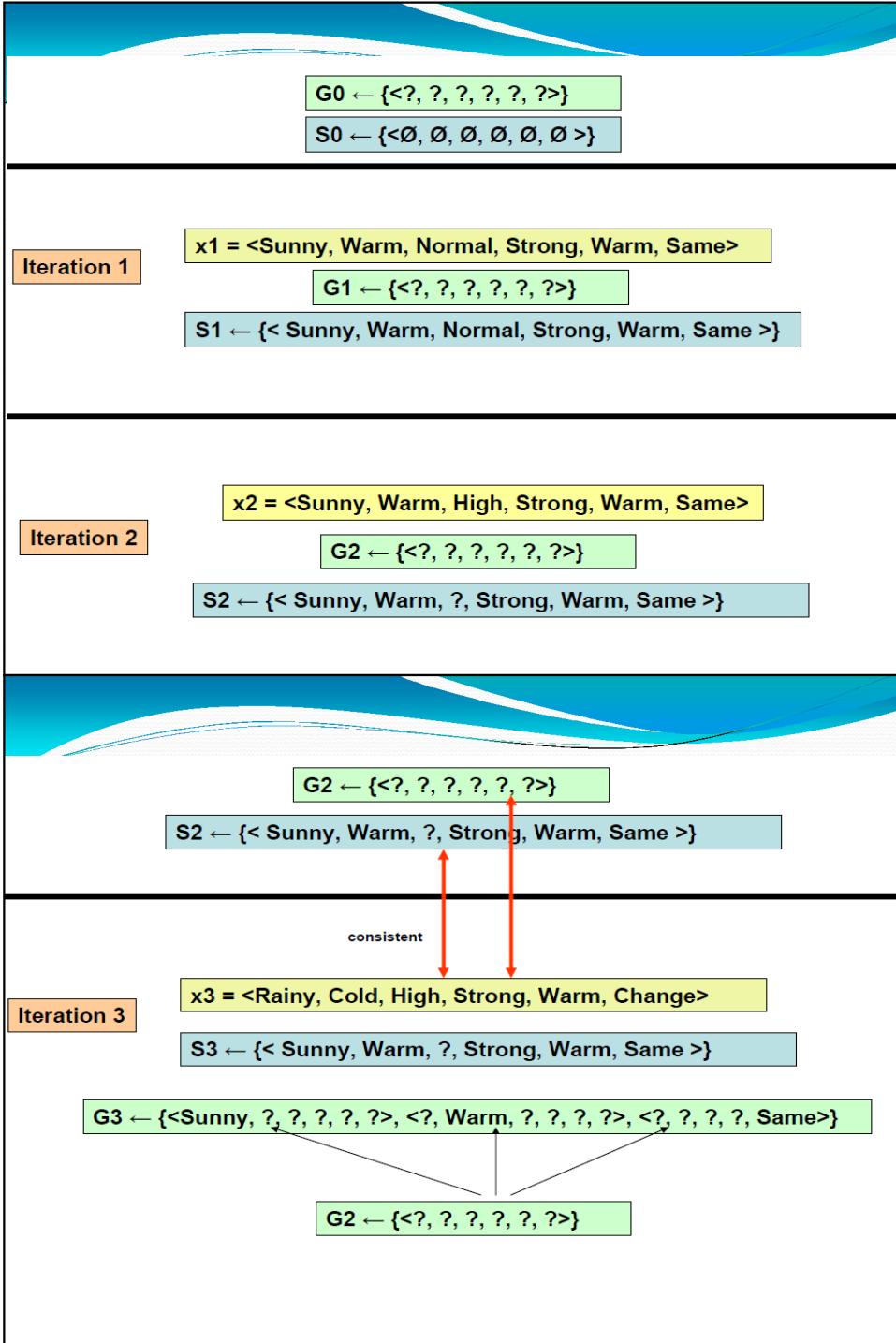
Example :

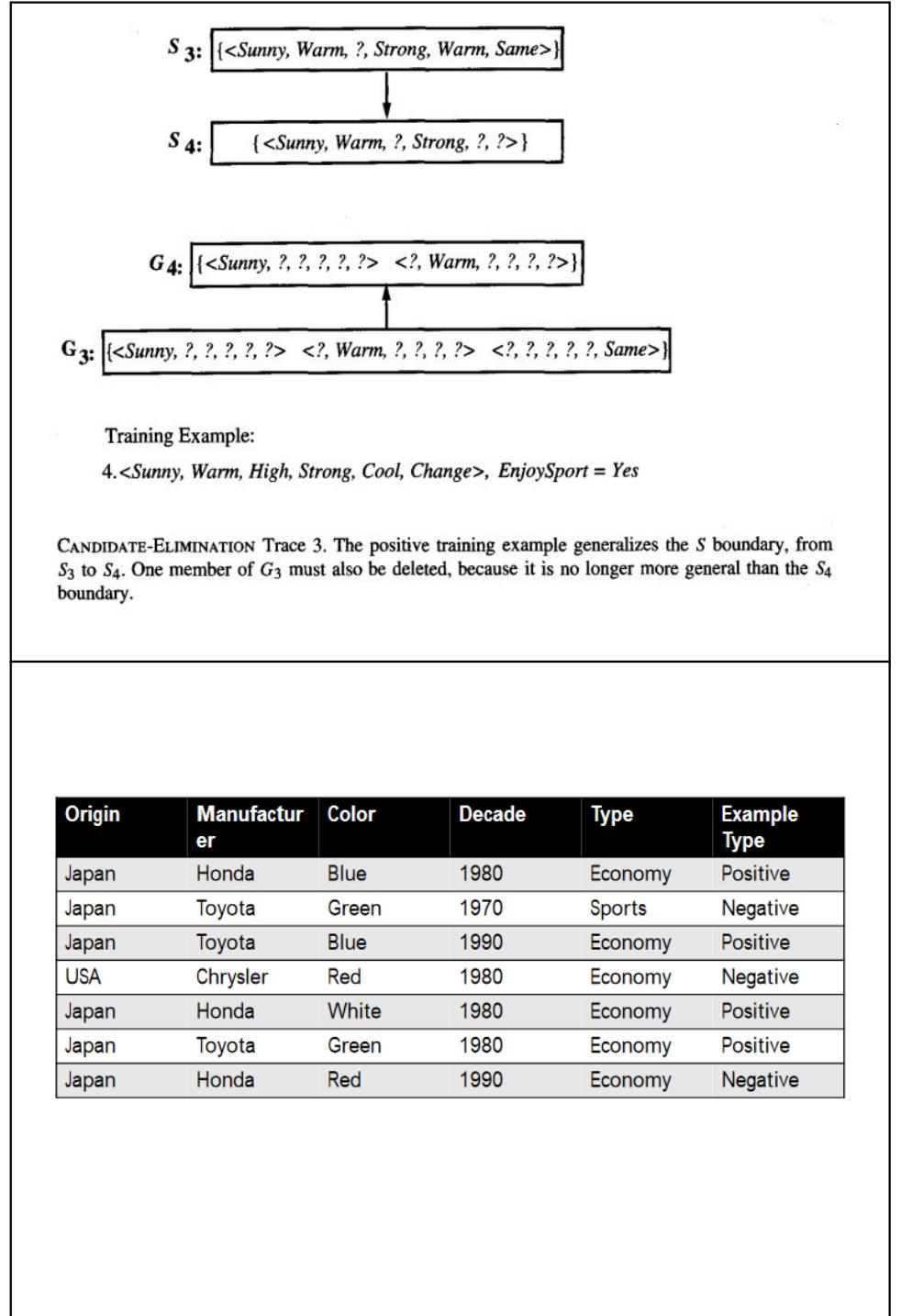
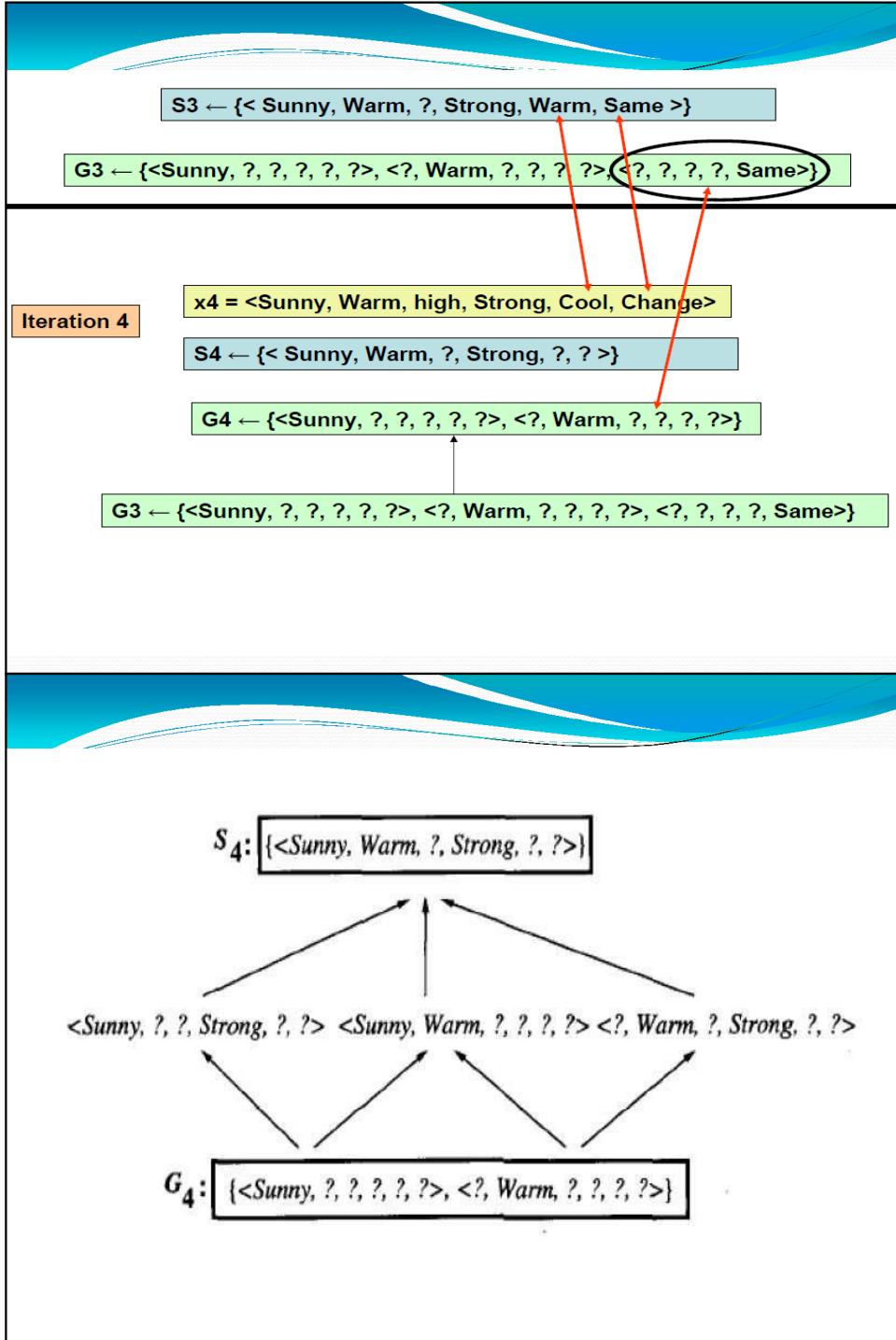
Example	Sky	AirTemp	Humidity	Wind	Water	Forecast	EnjoySport
1	Sunny	Warm	Normal	Strong	Warm	Same	Yes
2	Sunny	Warm	High	Strong	Warm	Same	Yes
3	Rainy	Cold	High	Strong	Warm	Change	No
4	Sunny	Warm	High	Strong	Cool	Change	Yes

$G_0 \leftarrow \{<?, ?, ?, ?, ?, ?>\}$

Initialization

$S_0 \leftarrow \{\emptyset, \emptyset, \emptyset, \emptyset, \emptyset, \emptyset\}$





What will Happen if the Training Contains errors ?

Example	Sky	AirTemp	Humidity	Wind	Water	Forecast	EnjoySport
1	Sunny	Warm	Normal	Strong	Warm	Same	Yes
2	Sunny	Warm	High	Strong	Warm	Same	No
3	Rainy	Cold	High	Strong	Warm	Change	No
4	Sunny	Warm	High	Strong	Cool	Change	Yes

$G2 \leftarrow \{?, ?, Normal, ?, ?, ?\}$

$S2 \leftarrow \{Sunny, Warm, Normal, Strong, Warm, Same\}$

consistent

$x3 = \{Rainy, Cold, High, Strong, Warm, Change\}$

$S3 \leftarrow \{Sunny, Warm, Normal, Strong, Warm, Same\}$

$G3 \leftarrow \{?, ?, Normal, ?, ?, ?\}$

$G0 \leftarrow \{?, ?, ?, ?, ?, ?\}$

$S0 \leftarrow \{\emptyset, \emptyset, \emptyset, \emptyset, \emptyset, \emptyset\}$

Iteration 1

$x1 = \{Sunny, Warm, Normal, Strong, Warm, Same\}$

$G1 \leftarrow \{?, ?, ?, ?, ?, ?\}$

$S1 \leftarrow \{Sunny, Warm, Normal, Strong, Warm, Same\}$

Iteration 2

$x2 = \{Sunny, Warm, High, Strong, Warm, Same\}$

$G2 \leftarrow \{?, ?, Normal, ?, ?, ?\}$

$S2 \leftarrow \{Sunny, Warm, Normal, Strong, Warm, Same\}$

Iteration 3

$S3 \leftarrow \{Sunny, Warm, Normal, Strong, Warm, Same\}$

$G3 \leftarrow \{?, ?, Normal, ?, ?, ?\}$

Iteration 4

$x4 = \{Sunny, Warm, High, Strong, Cool, Change\}$

$S4 \leftarrow \{\}$

$G4 \leftarrow \{\}$

Empty

$G3 \leftarrow \{?, ?, Normal, ?, ?, ?\}$

Concept Learning

- **Concept learning** is the task of automatically inferring the general definition of some concept, given examples labelled as members or non members of the concept.
 - Inferring a boolean-valued function from training examples of its input and output.

Concept Learning

- We incrementally learn general concepts from specific training examples.
- Each concept can be viewed as **describing some subset of objects or events defined over a larger set.**
- Concept : A boolean-valued function defined over a larger set.
 - E.g.: A function defined over all animals, whose value is true for birds and false for other animals.

Concept Learning

- **Concept learning** is the task of automatically inferring the general definition of some concept, given examples labelled as members or non members of the concept.
 - Inferring a boolean-valued function from training examples of its input and output.

Concept Learning - Task

- Example: Task of learning the target concept “**Days on which my friend Aldo enjoys his favorite water sport**”
Positive and negative training examples for the target concept “*EnjoySport*”

Sky	Temp	Humid	Wind	Water	Forecst	EnjoySpt
Sunny	Warm	Normal	Strong	Warm	Same	Yes
Sunny	Warm	High	Strong	Warm	Same	Yes
Rainy	Cold	High	Strong	Warm	Change	No
Sunny	Warm	High	Strong	Cool	Change	Yes
- The attribute ***EnjoySport*** indicates whether or not Aldo enjoys his favorite water sport on this day.
- **The task is to learn to predict the value of *EnjoySport* for an arbitrary day, based on the values of its other attributes.**

Hypothesis Representation

- Hypothesis : h , a conjunction of constraints on the instance attributes.
- Let each hypothesis be a vector of six constraints, specifying the values of the six attributes (**Sky**, **AirTemp**, **Humidity**, **Wind**, **Water**, and **Forecast**).
- Each constraint can be
 - A specific value (e.g., Water = Warm)
 - Don't care (e.g., Water = ?)
 - No value allowed (e.g., Water = " ϕ ")
- Example: Sky, AirTemp, Humidity, Wind, Water, Forecast
 $< \text{Sunny}, ?, \text{High}, ?, ?, ? >$

Hypothesis Representation

- The **most general hypothesis**-that every day is a positive example - is represented by
 $< ?, ?, ?, ?, ?, ? >$
- The **most specific possible hypothesis**-that no day is a positive example-is represented by
 $< \phi, \phi, \phi, \phi, \phi, \phi >$

Hypothesis Representation

- If some instance x satisfies all the constraints of hypothesis h , then h classifies x as a **positive example** ($h(x) = 1$).
- Example: The hypothesis that Aldo enjoys his favorite sport only on cold days with high humidity (independent of the values of the other attributes) is represented by the expression
 $< ?, \text{Cold}, \text{High}, ?, ?, ? >$

Hypothesis Representation

- In general, any concept learning task can be described by
 - the **set of instances** over which the target function is defined
 - the **target function**
 - the **set of candidate hypotheses** considered by the learner, and
 - the **set of available training examples**

Prototypical Concept Learning Task

- Given:

- Instances X : Possible days, each described by the attributes *Sky*, *AirTemp*, *Humidity*, *Wind*, *Water*, *Forecast*
- Target function c : $\text{EnjoySport} : X \rightarrow \{0, 1\}$
- Hypotheses H : Conjunctions of literals. E.g.
 $\langle ?, \text{Cold}, \text{High}, ?, ?, ? \rangle$.
- Training examples D : Positive and negative examples of the target function
 $\langle x_1, c(x_1) \rangle, \dots, \langle x_m, c(x_m) \rangle$

- Determine: A hypothesis h in H such that $h(x) = c(x)$ for all x in D .

Concept Learning Task- EnjoySport Prototype

- The set of items over which the concept is defined is called the set of instances, which we denote by X .
- Example: X is the set of all possible days, each represented by the attributes Sky, AirTemp, Humidity, Wind, Water, and Forecast.
- The concept or function to be learned is called the target concept, which we denote by c .
 - In general, c can be any boolean-valued function defined over the instances X ; that is, $c : X \rightarrow \{0, 1\}$
- Example: The target concept corresponds to the value of the attribute EnjoySport
 - (i.e., $c(x) = 1$ if EnjoySport = Yes, and $c(x) = 0$ if EnjoySport = No)

Concept Learning Task- EnjoySport Prototype

- Given:

- Instances X : Possible days, each described by the attributes
 - *Sky* (with possible values *Sunny*, *Cloudy*, and *Rainy*),
 - *AirTemp* (with values *Warm* and *Cold*),
 - *Humidity* (with values *Normal* and *High*),
 - *Wind* (with values *Strong* and *Weak*),
 - *Water* (with values *Warm* and *Cool*), and
 - *Forecast* (with values *Same* and *Change*).
- Hypotheses H : Each hypothesis is described by a conjunction of constraints on the attributes *Sky*, *AirTemp*, *Humidity*, *Wind*, *Water*, and *Forecast*. The constraints may be “?” (any value is acceptable), “ \emptyset ” (no value is acceptable), or a specific value.
- Target concept c : $\text{EnjoySport} : X \rightarrow \{0, 1\}$
- Training examples D : Positive and negative examples of the target function

- Determine:

- A hypothesis h in H such that $h(x) = c(x)$ for all x in X .

Concept Learning Task- EnjoySport Prototype

- When learning the target concept, the learner is presented a set of training examples, each consisting of an instance x from X , along with its target concept value $c(x)$
- Instances for which $c(x) = 1$ are called positive examples, or members of the target concept.
- Instances for which $c(x) = 0$ are called negative examples, or non members of the target concept.
- Ordered pair $\langle x, c(x) \rangle$ describes the training example consisting of the instance x and its target concept value $c(x)$.
- D denotes the set of available training examples.

Concept Learning Task- EnjoySport Prototype

- Given a set of training examples of the target concept c , **the problem faced by the learner is to hypothesize, or estimate, c .**
- H denotes the set of **all possible hypotheses** that the learner may consider regarding the identity of the target concept.
- In general, **each hypothesis h in H represents a boolean-valued function defined over X** ; that is, $h : X \rightarrow \{0, 1\}$.
- The goal of the learner is to find a hypothesis h such that $h(x) = c(x)$ for all x in X .

CONCEPT LEARNING AS SEARCH

- Concept Learning : Viewed as the **task of searching through a large space of hypotheses** implicitly defined by the hypothesis representation.
- Goal of search:** Find the hypothesis that best fits the training examples.
 - By selecting a hypothesis representation, the designer of the learning algorithm implicitly defines the space of all hypotheses that the program can ever represent and therefore can ever learn.
- Example: Consider the instance X and hypotheses H in the **EnjoySport** learning task.
 - Given that the attribute **Sky** has three possible values, and that **AirTemp, Humidity, Wind, Water, and Forecast** each have two possible values, the **instance space X contains exactly $3 \cdot 2 \cdot 2 \cdot 2 \cdot 2 = 96$ distinct instances**.
 - Further, there are **5.4.4.4.4.4 = 5120** syntactically distinct hypotheses within H .

The Inductive Learning Hypothesis

- Learning Task :**To determine a hypothesis h identical to the target concept c over the entire set of instances X**
 - The only information available about c is its value over the training examples.**
- Inductive learning algorithms guarantee that the output hypothesis fits the target concept over the training data.
- Assumption : The best hypothesis regarding unseen instances is the hypothesis that best fits the observed training data.
- The inductive learning hypothesis: **Any hypothesis found to approximate the target function well over a sufficiently large set of training examples will also approximate the target function well over other unobserved examples.**

CONCEPT LEARNING AS SEARCH

- Every hypothesis containing one or more ϕ symbols represents the empty set of instances (classified as negative)
- Therefore, the number of semantically distinct hypotheses is only **1+ (4.3.3.3.3.3)**
- EnjoySport example is a very simple learning task, with a relatively small, **finite hypothesis space**.
- Most practical learning tasks involve much larger, sometimes **infinite hypothesis spaces**.

FIND-S: FINDING A MAXIMALLY specific HYPOTHESIS

- Use the *more_general_than* partial ordering to organize the search for a hypothesis consistent with the observed training examples
 - Begin with the most specific possible hypothesis in H , then generalize this hypothesis each time it fails to cover an observed positive training example.
 - A hypothesis "covers" a positive example if it correctly classifies the example as positive
- FIND-S algorithm
 1. Initialize h to the most specific hypothesis in H
 2. For each positive training instance x
 - For each attribute constraint a , in h
 - If the constraint a , is satisfied by x Then do nothing
 - Else replace a , in h by the next more general constraint that is satisfied by x
 3. Output hypothesis h

FIND-S: Example

Hence, $h \leftarrow (\text{Sunny}, \text{Warm}, \text{Normal}, \text{Strong}, \text{Warm}, \text{Same})$

- This h is still very specific; it asserts that all instances are negative except for the single positive training example observed.
- Next, the second training example (positive example) forces the algorithm to further generalize h , this time substituting a "?" in place of any attribute value in h that is not satisfied by the new example.
- The refined hypothesis in this case is
 $h \leftarrow (\text{Sunny}, \text{Warm}, ?, \text{Strong}, \text{Warm}, \text{Same})$
- The third training example (a negative example) the algorithm makes no change to h .
- FIND-S algorithm simply ignores every negative example, and hence no revision in h , is needed.

Sky	Temp	Humid	Wind	Water	Forecst	EnjoySpt
Sunny	Warm	Normal	Strong	Warm	Same	Yes
Sunny	Warm	High	Strong	Warm	Same	Yes
Rainy	Cold	High	Strong	Warm	Change	No
Sunny	Warm	High	Strong	Cool	Change	Yes

FIND-S: Example

- Illustration: *EnjoySport* task
 - The first step of FIND-S is to initialize h to the most specific hypothesis in H

$$h \leftarrow <\phi, \phi, \phi, \phi, \phi, \phi>$$
 - Observing the first training example in *EnjoySport*, which is a positive example, it is clear that the hypothesis is too specific.
 - None of the " ϕ " constraints in h are satisfied by this example, so each is replaced by the next more general constraint that fits the example; namely, the attribute values for this training example.

Sky	Temp	Humid	Wind	Water	Forecst	EnjoySpt
Sunny	Warm	Normal	Strong	Warm	Same	Yes
Sunny	Warm	High	Strong	Warm	Same	Yes
Rainy	Cold	High	Strong	Warm	Change	No
Sunny	Warm	High	Strong	Cool	Change	Yes

FIND-S: Example

- The fourth (positive) example leads to a further generalization of h .

$h \leftarrow (\text{Sunny}, \text{Warm}, ?, \text{Strong}, ?, ?)$

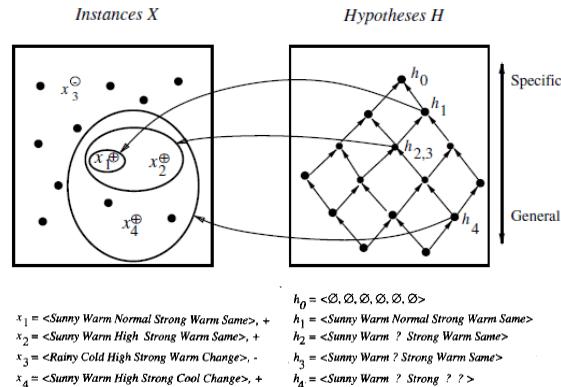
- The search moves from hypothesis to hypothesis, searching from the most specific to progressively more general hypotheses.

Sky	Temp	Humid	Wind	Water	Forecst	EnjoySpt
Sunny	Warm	Normal	Strong	Warm	Same	Yes
Sunny	Warm	High	Strong	Warm	Same	Yes
Rainy	Cold	High	Strong	Warm	Change	No
Sunny	Warm	High	Strong	Cool	Change	Yes

FIND-S: Example

Search in terms of the instance and hypothesis spaces:

- At each step, the hypothesis is generalized only as far as necessary to cover the new positive example.
- At each stage the hypothesis is the most specific hypothesis consistent with the training examples observed up to this point (hence the name FIND-S).



Find-S Example

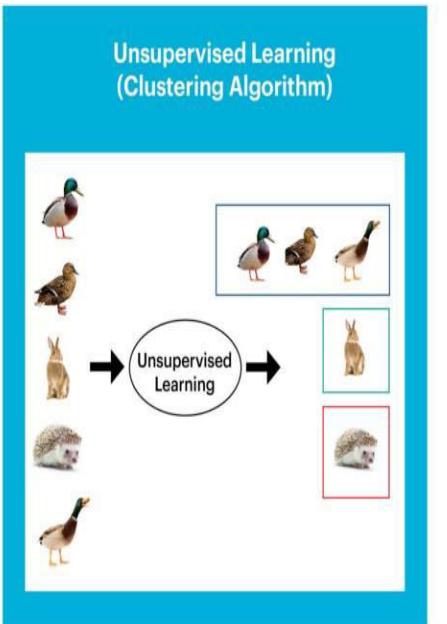
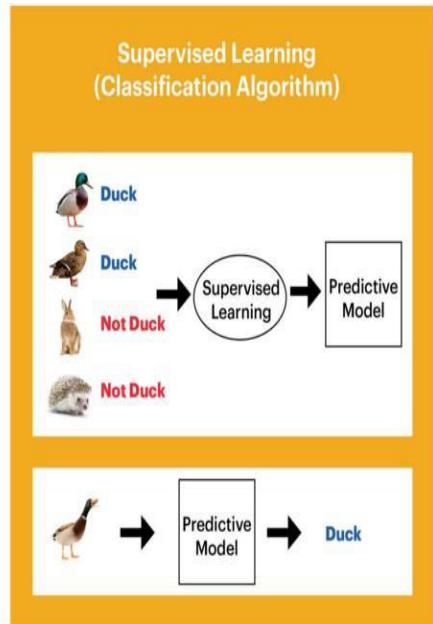
Eyes	Nose	Head	Fcolor	Hair?	Smile?
Round	Triangle	Round	Purple	Yes	Yes
Square	Square	Square	Green	Yes	No
Square	Triangle	Round	Yellow	Yes	Yes
Round	Triangle	Round	Green	No	No
Square	Square	Round	Yellow	Yes	Yes

FIND-S: Example

- The key property of the FIND-S algorithm is that for hypothesis spaces described by conjunctions of attribute constraints (such as H for the **EnjoySport** task),
 - FIND-S is guaranteed to output the most specific hypothesis within H that is consistent with the positive training examples.
- Its final hypothesis will also be consistent with the negative examples provided the correct target concept is contained in H , and provided the training examples are correct.

Supervised vs. Unsupervised Learning

- **Supervised learning (classification)**
 - Supervision: The training data (observations, measurements, etc.) are accompanied by **labels** indicating the class of the observations
 - New data is classified based on the training set
- **Unsupervised learning (clustering)**
 - The class labels of training data is unknown
 - Given a set of measurements, observations, etc. with the aim of establishing the existence of classes or clusters in the data



CLASSIFICATION

Classification vs. Prediction

- **Classification**
 - predicts categorical class labels (**discrete** or nominal)
 - classifies data (constructs a model) based on the training set and the values (**class labels**) in a classifying attribute and uses it in classifying new data
- **Prediction**
 - models **continuous-valued** functions, i.e., predicts unknown or missing values
- Typical applications
 - Credit approval
 - Target marketing
 - Medical diagnosis
 - Fraud detection

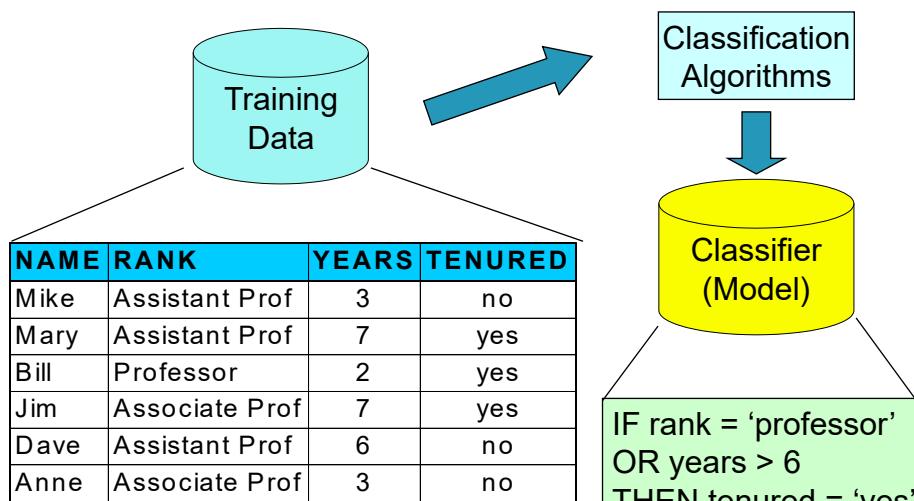
Classification: Definition

- Given a collection of records (**training set**)
 - Each record contains a set of **attributes**, one of the attributes is the **class**.
- Find a **model** for class attribute as a function of the values of other attributes.
- Goal: previously unseen records should be assigned a class as accurately as possible.
 - A **test set** is used to determine the accuracy of the model. Usually, the given data set is divided into training and test sets, with training set used to build the model and test set used to validate it.

Classification—A Two-Step Process

- **Model construction:** describing a set of predetermined classes
 - Each tuple/sample is assumed to belong to a predefined class, as determined by the **class label attribute**
 - The set of tuples used for model construction is **training set**
 - The model is represented as **classification rules, decision trees, or mathematical formulae**

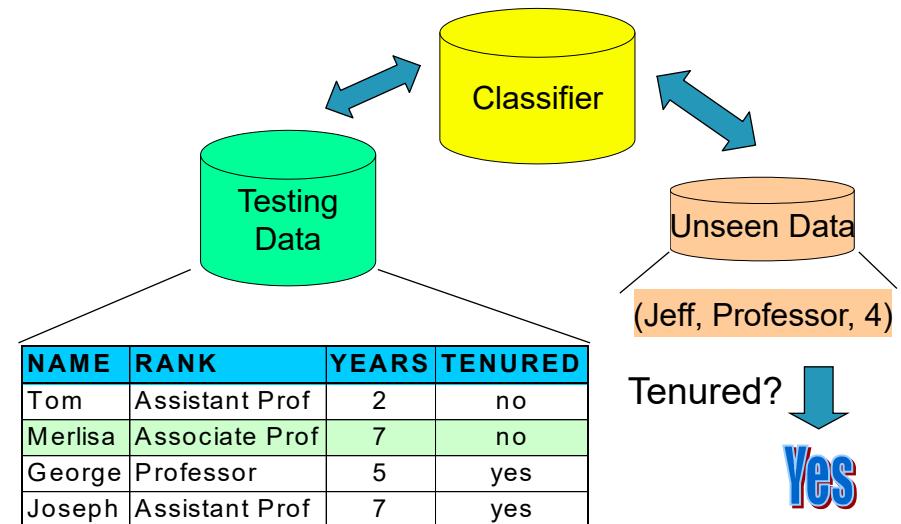
Classification Process (1): Model Construction



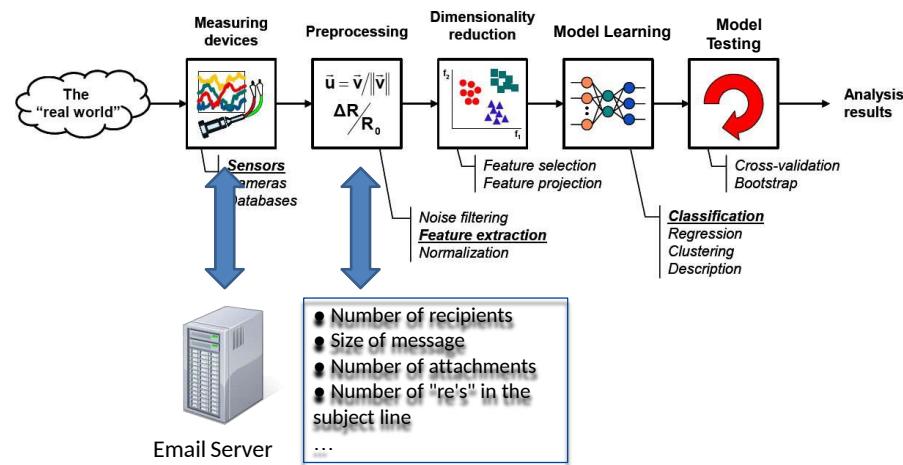
Classification—A Two-Step Process

- **Model usage:** for classifying future or unknown objects
 - Estimate accuracy of the model
 - The known label of test sample is compared with the classified result from the model
 - Accuracy rate is the percentage of test set samples that are correctly classified by the model
 - Test set is independent of training set, otherwise over-fitting will occur
 - If the accuracy is acceptable, use the model to **classify data** tuples whose class labels are not known

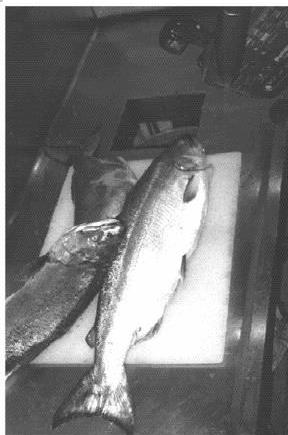
Classification Process (2): Use the Model in Prediction



The Learning Process in spam mail Example



An Example (continued)



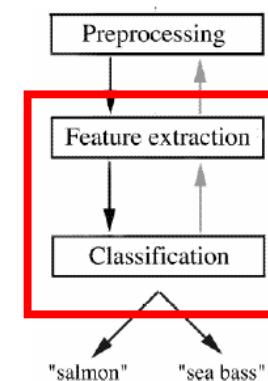
Features/attributes:

- Length
- Lightness
- Width
- Position of mouth

An Example

- A fish-packing plant wants to automate the process of sorting incoming fish according to species
- As a pilot project, it is decided to try to separate **sea bass** from **salmon** using optical sensing

An Example (continued)



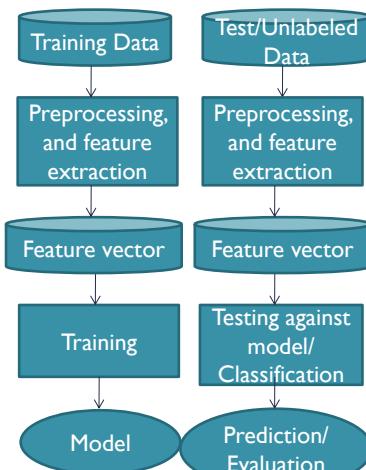
- **Preprocessing:** Images of different fishes are isolated from one another and from background;
- **Feature extraction:** The information of a single fish is then sent to a feature extractor, that measure certain "features" or "properties";
- **Classification:** The values of these features are passed to a classifier that evaluates the evidence presented, and build a model to discriminate between the two species

An Example (continued)

- Domain knowledge:
 - A sea bass is generally longer than a salmon
- Related feature: (or attribute)
 - Length
- Training the classifier:
 - Some examples are provided to the classifier in this form: <fish_length, fish_name>
 - These examples are called training examples
 - The classifier *learns* itself from the training examples, how to distinguish Salmon from Bass based on the fish_length

An Example (continued)

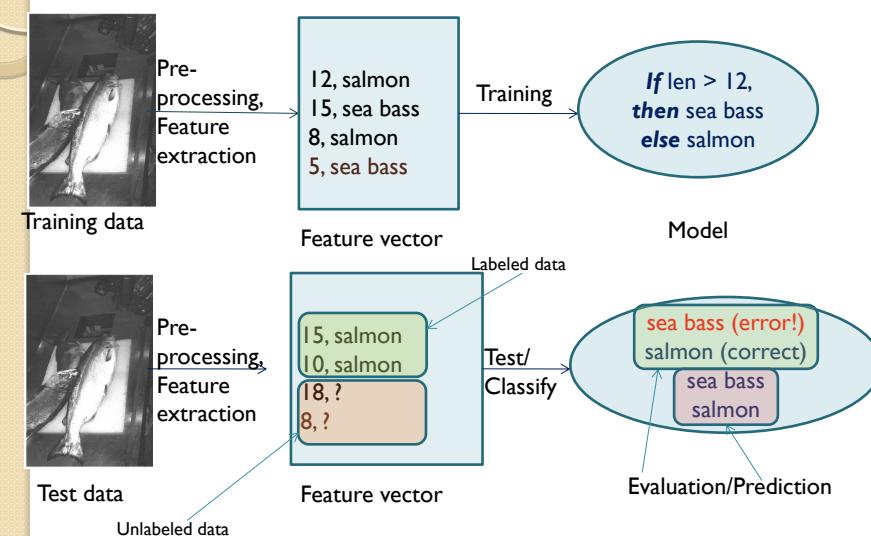
- So the overall classification process goes like this →



An Example (continued)

- Classification model (hypothesis):
 - The classifier generates a model from the training data to classify future examples (test examples)
 - An example of the model is a rule like this:
 - If $\text{Length} \geq l^*$ then sea bass otherwise salmon
 - Here the value of l^* determined by the classifier
- Testing the model
 - Once we get a model out of the classifier, we may use the classifier to test future examples
 - The test data is provided in the form <fish_length>
 - The classifier outputs <fish_type> by checking fish_length against the model

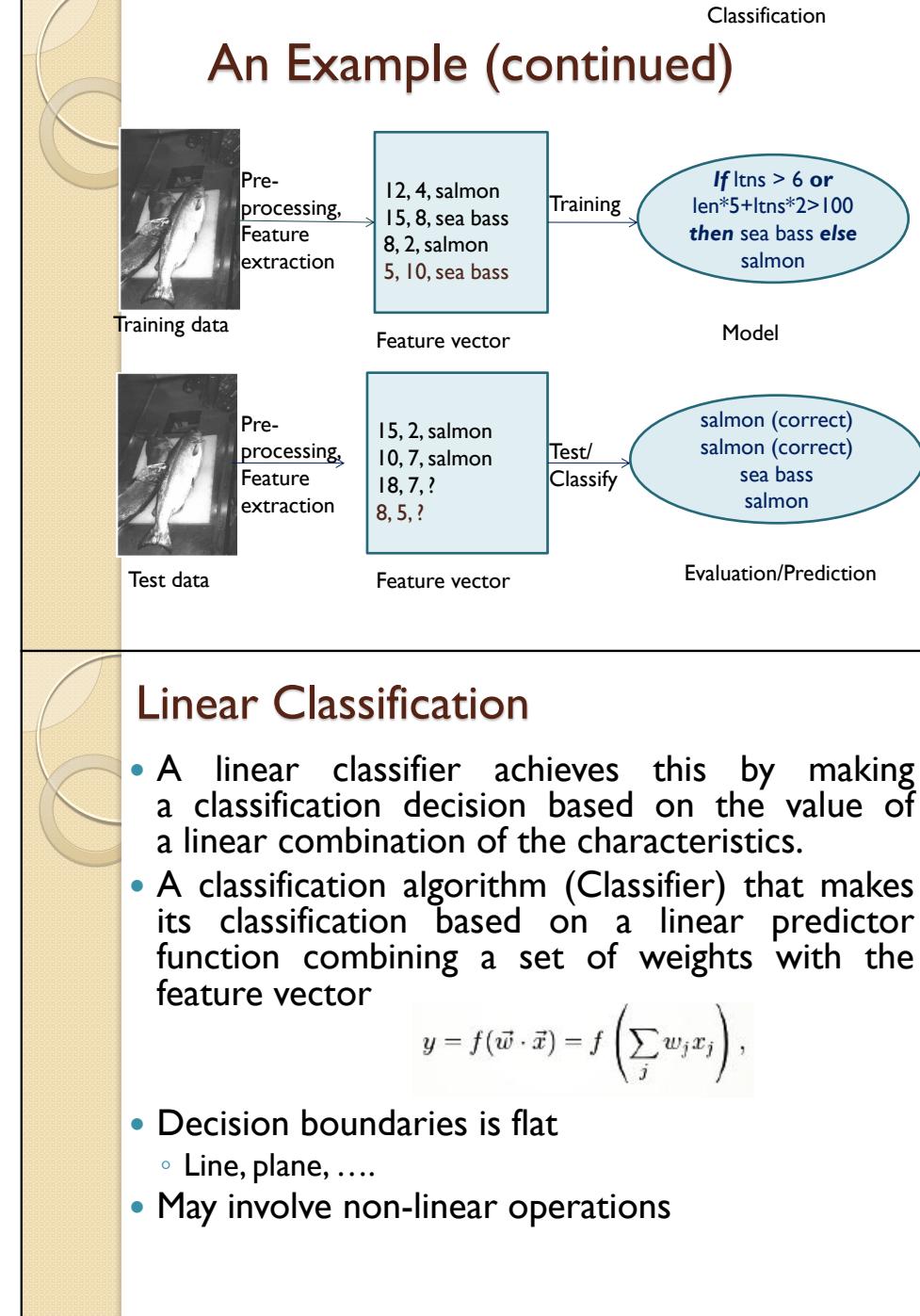
An Example (continued)



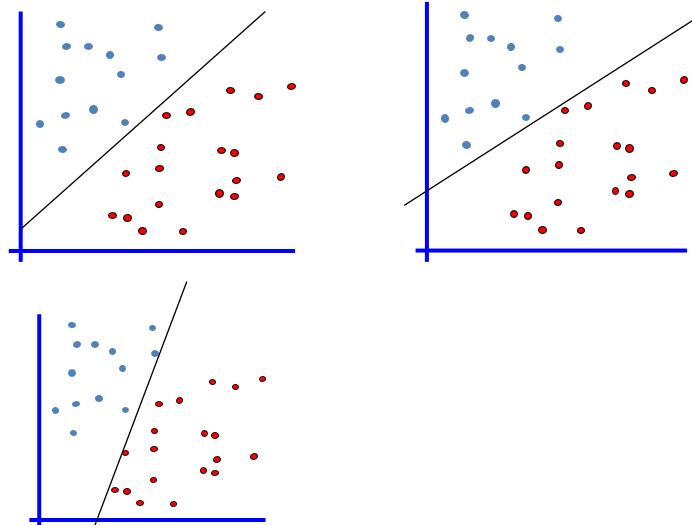
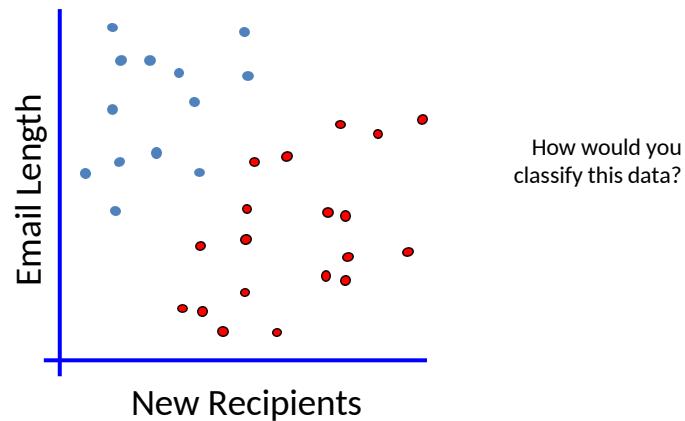
An Example (continued)

- Why error?
 - Insufficient training data
 - Too few features
 - Too many/irrelevant features
 - Overfitting / specialization

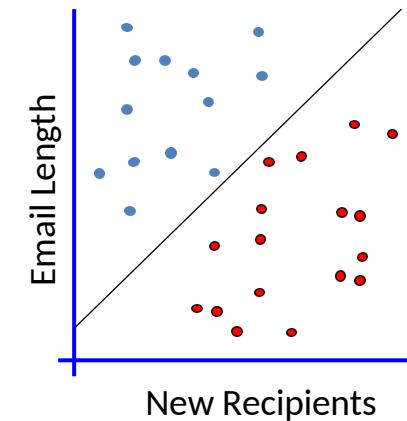
Linear, Non-linear, Multi-class
and
Multi-label classification



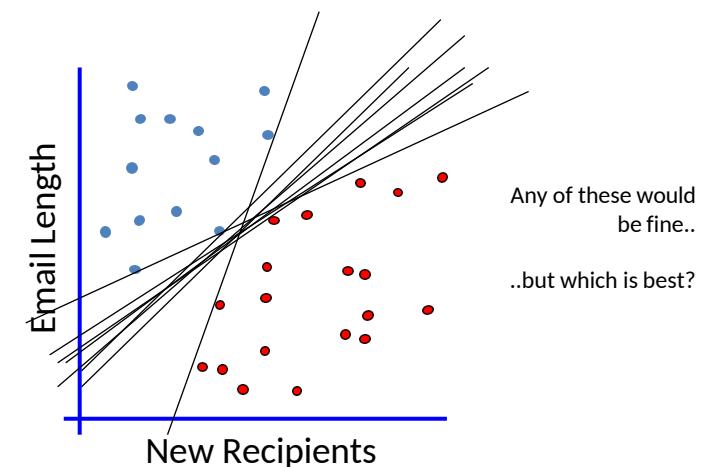
Linear Classifiers



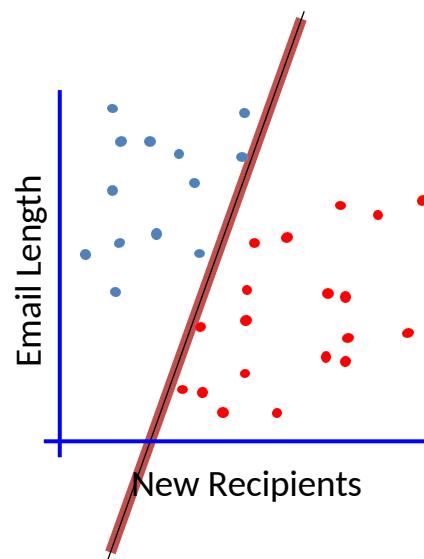
Linear Classifiers



Linear Classifiers



Classifier Margin



Define the **margin** of a linear classifier as the width that the boundary could be increased by before hitting a datapoint.

Acceptance at a University



Acceptance at a University



Test



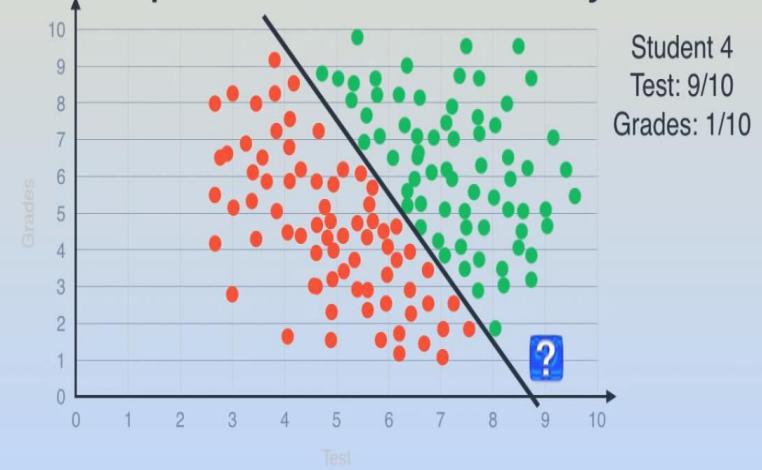
Grades

Student 1
Test: 9/10
Grades: 8/10

Student 2
Test: 3/10
Grades: 4/10

Student 3
Test: 7/10
Grades: 6/10

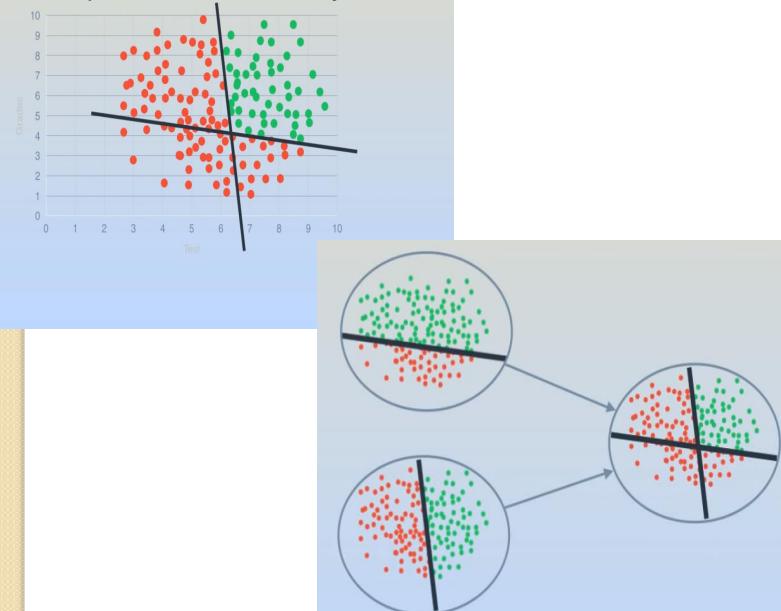
Acceptance at a University



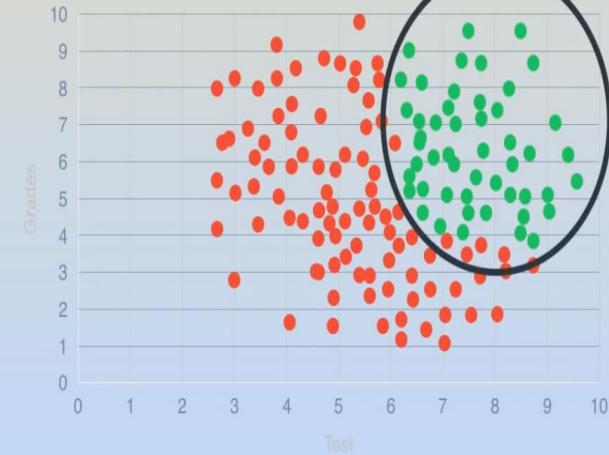
Acceptance at a University



Acceptance at a University



Acceptance at a University

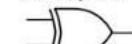


XOR & XNOR Gate: Truth Table & Symbol

XOR / XNOR Tables and Symbols

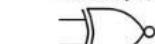
X	Y	$X \oplus Y$
0	0	0
0	1	1
1	0	1
1	1	0

XOR Symbol

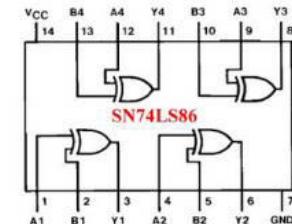


X	Y	$\overline{X} \oplus \overline{Y}$
0	0	1
0	1	0
1	0	0
1	1	1

XNOR Symbol



The XNOR is also denoted as equivalence

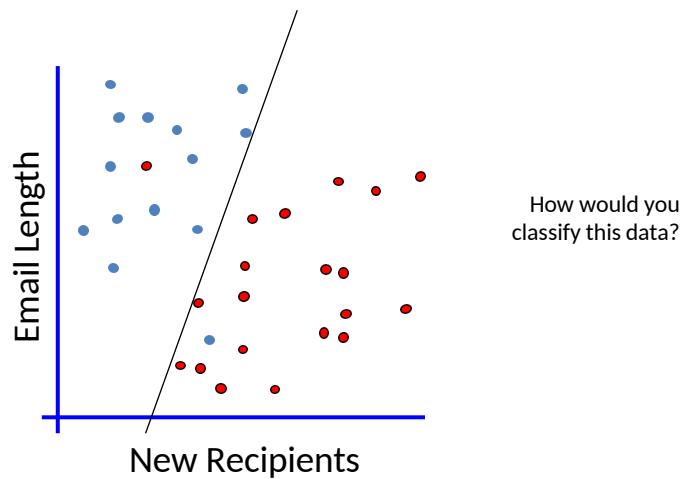


Inputs		Output
A	B	Y
L	L	L
L	H	H
H	L	H
H	H	L

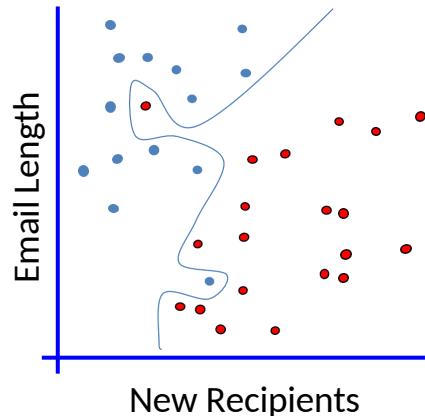


Electrical 4 U

No Linear Classifier can cover all instances



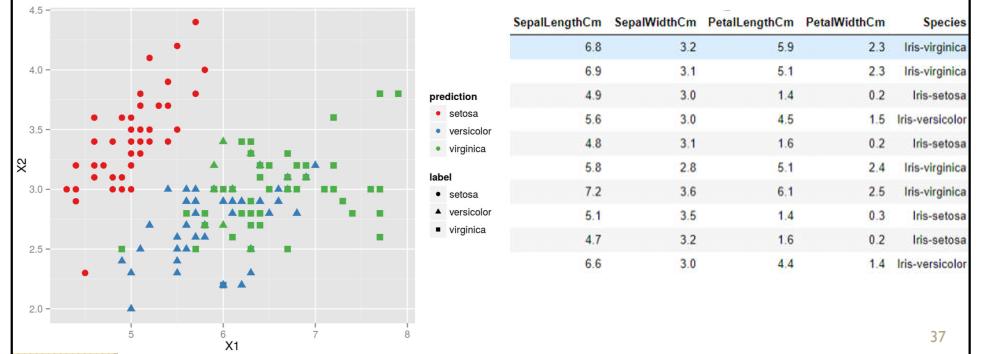
No Linear Classifier can cover all instances



- Ideally, the best decision boundary should be the one which provides an optimal performance such as in the following figure

What is multiclass

- Output $\in \{1, 2, 3, \dots, K\}$
 - In some cases, output space can be very large (i.e., K is very large)
- Each input belongs to exactly one class (c.f. in multilabel, input belongs to many classes)



Multi-Classes Classification

- Multi-class classification is simply classifying objects into any one of multiple categories. Such as classifying just into either a dog or cat from the dataset.
- 1. When there are more than two categories in which the images can be classified, and
- 2. An image does not belong to more than one class
-
- If both of the above conditions are satisfied, it is referred to as a multi-class image classification problem



Multi-label classification

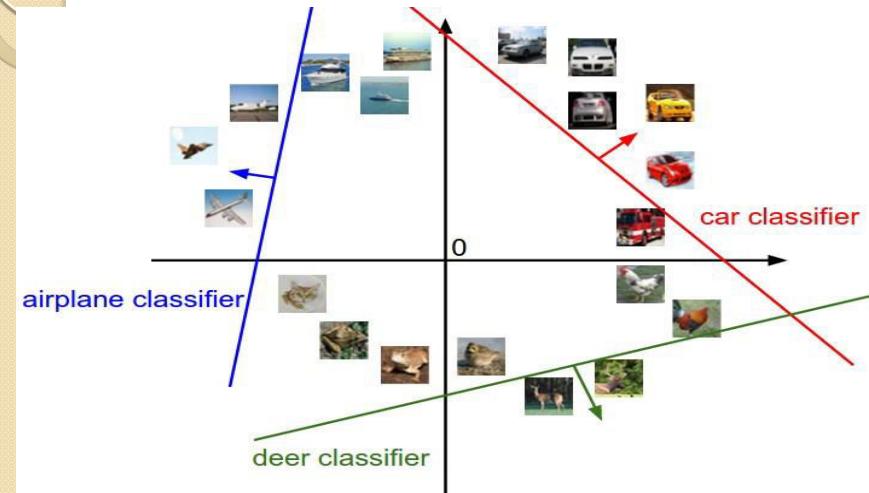
- When we can classify an image into more than one class (as in the image beside), it is known as a multi-label image classification problem.



- **Multi-label classification** is a type of **classification** in which an object can be categorized into more than one class.

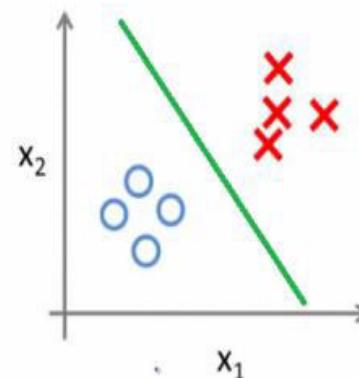
- For example, In the image dataset, we will **classify** a picture as the **image** of a dog or cat and also **classify** the same **image** based on the breed of the dog or cat

These are all labels of the given images. **Each image here belongs to more than one class and hence it is a multi-label image classification problem.**

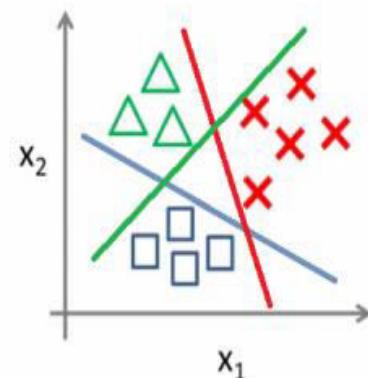


Binary Vs Multi-class

Binary classification:



Multi-class classification:



Three Type of Classification Tasks

Binary Classification



- Spam
- Not spam

Multiclass Classification



- Dog
- Cat
- Horse
- Fish
- Bird
- ...

YAHOO!
JAPAN

Multi-label Classification



- Dog
- Cat
- Horse
- Fish
- Bird
- ...

Multi class Vs multi label classification

Table 1

x	y
X_1	t_1
X_2	t_2
X_3	t_1
X_4	t_2
X_5	t_1

Binary Classification

Table 2

x	y
X_1	t_2
X_2	t_3
X_3	t_4
X_4	t_1
X_5	t_3

Multi-class Classification

Table 3

x	y
X_1	$[t_2, t_5]$
X_2	$[t_1, t_2, t_3, t_4]$
X_3	$[t_3]$
X_4	$[t_2, t_4]$
X_5	$[t_1, t_3, t_4]$

Multi-label Classification

Decision Tree Induction (ID3/C4.5)

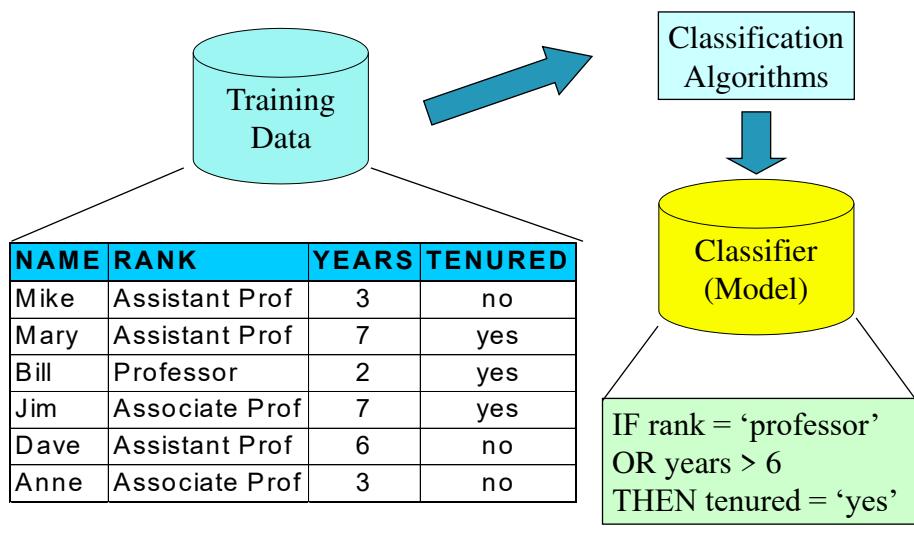
Supervised vs. Unsupervised Learning

- Supervised learning (classification)
 - Supervision: The training data (observations, measurements, etc.) are accompanied by **labels** indicating the class of the observations
 - New data is classified based on the training set
- Unsupervised learning (clustering)
 - The class labels of training data is unknown
 - Given a set of measurements, observations, etc. with the aim of establishing the existence of classes or clusters in the data

Prediction Problems: Classification vs. Numeric Prediction

- **Classification**
 - predicts categorical class labels (discrete or nominal)
 - classifies data (constructs a model) based on the training set and the values (**class labels**) in a classifying attribute and uses it in classifying new data
- **Numeric Prediction**
 - models continuous-valued functions, i.e., predicts unknown or missing values
- **Typical applications**
 - Credit/loan approval:
 - Medical diagnosis: if a tumor is cancerous or benign
 - Fraud detection: if a transaction is fraudulent
 - Web page categorization: which category it is

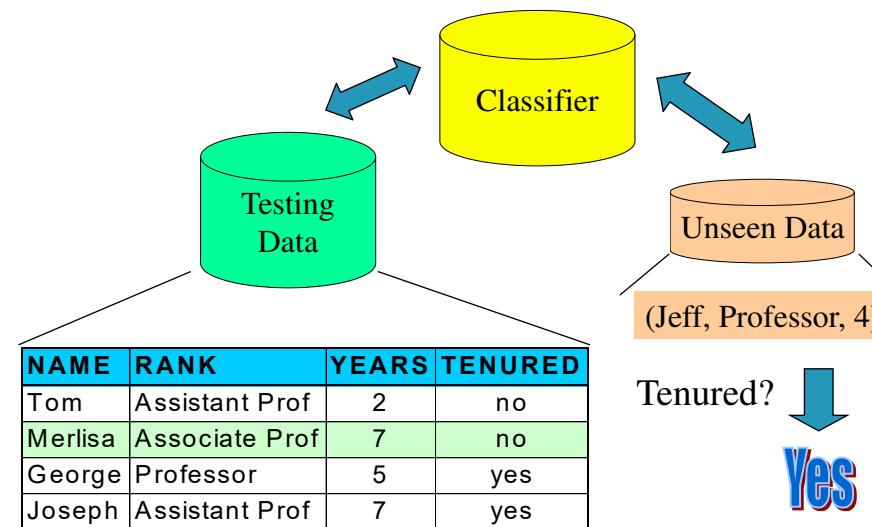
Process (1): Model Construction



Classification—A Two-Step Process

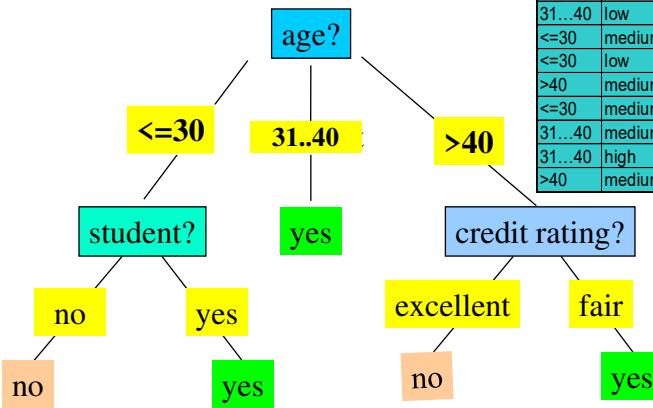
- **Model construction:** describing a set of predetermined classes
 - Each tuple/sample is assumed to belong to a predefined class, as determined by the **class label attribute**
 - The set of tuples used for model construction is **training set**
 - The model is represented as classification rules, decision trees, or mathematical formulae
- **Model usage:** for classifying future or unknown objects
 - **Estimate accuracy** of the model
 - The known label of test sample is compared with the classified result from the model
 - **Accuracy** rate is the percentage of test set samples that are correctly classified by the model
 - **Test set** is independent of training set (otherwise overfitting)
 - If the accuracy is acceptable, use the model to **classify new data**
 - Note: If the test set is used to select models, it is called **validation (test) set**

Process (2): Using the Model in Prediction



Decision Tree Induction: An Example

- Training data set: Buys_computer
- Resulting tree:



age	income	student	credit rating	buys computer
<=30	high	no	fair	no
<=30	high	no	excellent	no
31..40	high	no	fair	yes
>40	medium	no	fair	yes
>40	low	yes	fair	yes
>40	low	yes	excellent	no
31..40	low	yes	excellent	yes
<=30	medium	no	fair	no
<=30	low	yes	fair	yes
>40	medium	yes	fair	yes
<=30	medium	yes	excellent	yes
31..40	medium	no	excellent	yes
31..40	high	yes	fair	yes
>40	medium	no	excellent	no

Attribute Selection Measure: Information Gain (ID3/C4.5)

- Select the attribute with the highest information gain
- Let p_i be the probability that an arbitrary tuple in D belongs to class C_i , estimated by $|C_{i,D}|/|D|$
- Expected information** (entropy) needed to classify a tuple in D:

$$Info(D) = -\sum_{i=1}^m p_i \log_2(p_i)$$
- Information** needed (after using A to split D into v partitions) to classify D:

$$Info_A(D) = \sum_{j=1}^v \frac{|D_j|}{|D|} \times Info(D_j)$$
- Information gained** by branching on attribute A

$$Gain(A) = Info(D) - Info_A(D)$$

Algorithm for Decision Tree Induction

- Basic algorithm (a greedy algorithm)
 - Tree is constructed in a **top-down recursive divide-and-conquer manner**
 - At start, all the training examples are at the root
 - Attributes are categorical (if continuous-valued, they are discretized in advance)
 - Examples are partitioned recursively based on selected attributes
 - Test attributes are selected on the basis of a heuristic or statistical measure (e.g., **information gain**)
- Conditions for stopping partitioning
 - All samples for a given node belong to the same class
 - There are no remaining attributes for further partitioning – **majority voting** is employed for classifying the leaf
 - There are no samples left

Attribute Selection: Information Gain

- Class P: buys_computer = "yes"
- Class N: buys_computer = "no"

$$Info(D) = -\sum_{i=1}^m p_i \log_2(p_i)$$

$$Info(D) = I(9,5) = -\frac{9}{14} \log_2(\frac{9}{14}) - \frac{5}{14} \log_2(\frac{5}{14}) = 0.940$$

age	p_i	n_i	$I(p_i, n_i)$
<=30	2	3	0.971
31..40	4	0	0
>40	3	2	0.971

age	income	student	credit rating	buys computer
<=30	high	no	fair	no
<=30	high	no	excellent	no
31..40	high	no	fair	yes
>40	medium	no	fair	yes
>40	low	yes	fair	yes
>40	low	yes	excellent	no
31..40	low	yes	excellent	yes
<=30	medium	no	fair	no
<=30	low	yes	fair	yes
>40	medium	yes	fair	yes
<=30	medium	yes	excellent	yes
31..40	medium	no	excellent	yes
31..40	high	yes	fair	yes
>40	medium	no	excellent	no

age	p_i	n_i	$I(p_i, n_i)$
<=30	2	3	0.971
31...40	4	0	0
>40	3	2	0.971

age	income	student	credit rating	buys computer
<=30	high	no	fair	no
<=30	high	no	excellent	no
31...40	high	no	fair	yes
>40	medium	no	fair	yes
>40	low	yes	fair	yes
>40	low	yes	excellent	no
31...40	low	yes	excellent	yes
<=30	medium	no	fair	no
<=30	low	yes	fair	yes
>40	medium	yes	fair	yes
<=30	medium	yes	excellent	yes
31...40	medium	no	excellent	yes
31...40	high	yes	fair	yes
>40	medium	no	excellent	no

$$\text{Info}_{age}(D) = \frac{5}{14} I(2,3) + \frac{4}{14} I(4,0) + \frac{5}{14} I(3,2) = 0.694$$

$$\text{Info}_{age}(D) = \frac{5}{14} \times \left(-\frac{2}{5} \log_2 \frac{2}{5} - \frac{3}{5} \log_2 \frac{3}{5} \right)$$

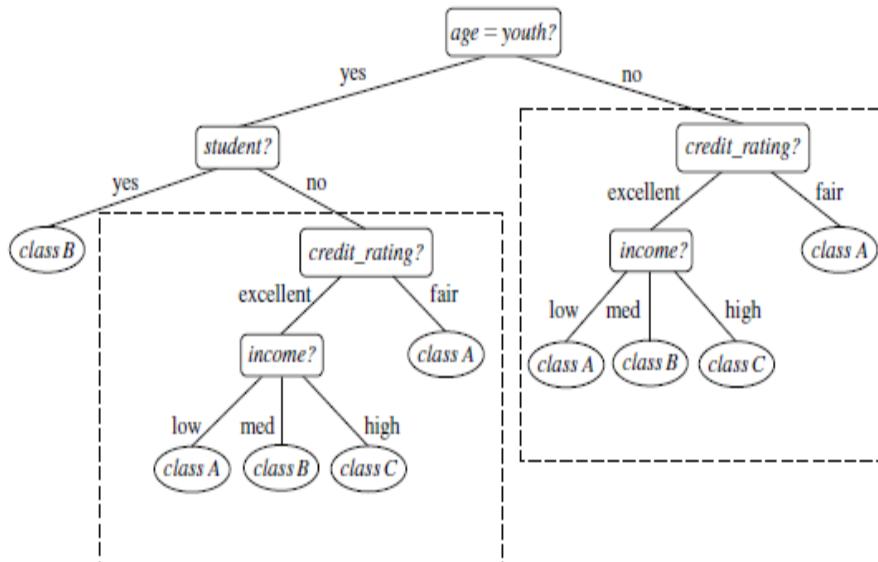
$\frac{5}{14} I(2,3)$ means "age <=30" has 5 out of 14 samples, with 2 yes'es and 3 no's. Hence

$$\text{Gain}(age) = \text{Info}(D) - \text{Info}_{age}(D) = 0.246$$

$$\text{Gain}(age) = \text{Info}(D) - \text{Info}_{age}(D) = 0.940 - 0.694 = 0.246 \text{ bits.}$$

$$+ \frac{4}{14} \times \left(-\frac{4}{4} \log_2 \frac{4}{4} \right) + \frac{5}{14} \times \left(-\frac{3}{5} \log_2 \frac{3}{5} - \frac{2}{5} \log_2 \frac{2}{5} \right)$$

$$= 0.694 \text{ bits.}$$



$$\text{Gain}(income) = 0.029$$

$$\text{Gain}(student) = 0.151$$

$$\text{Gain}(credit_rating) = 0.048$$

Create a decision tree for the data values given in the table using Entropy (Information Gain) based attribute finding.

Refund	Marital Status	Taxable Income	Cheat
Yes	Single	125K	No
No	Married	100K	No
No	Single	70K	No
Yes	Married	120K	No
No	Divorced	95K	Yes
No	Married	60K	No
Yes	Divorced	220K	No
No	Single	85K	Yes
No	Married	75K	No
No	Single	90K	Yes

Evaluating Classifier Accuracy: Bootstrap

- **Bootstrap**
 - Works well with small data sets
 - Samples the given training tuples uniformly *with replacement*
 - i.e., each time a tuple is selected, it is equally likely to be selected again and re-added to the training set
- Several bootstrap methods, and a common one is **.632 bootstrap**
 - A data set with d tuples is sampled d times, with replacement, resulting in a training set of d samples. The data tuples that did not make it into the training set end up forming the test set. About 63.2% of the original data end up in the bootstrap, and the remaining 36.8% form the test set (since $(1 - 1/d)^d \approx e^{-1} = 0.368$)
 - Repeat the sampling procedure k times, overall accuracy of the model:

$$Acc(M) = \frac{1}{k} \sum_{i=1}^k (0.632 \times Acc(M_i)_{test_set} + 0.368 \times Acc(M_i)_{train_set})$$

23

Estimating Confidence Intervals: Null Hypothesis

- Perform 10-fold cross-validation
- Assume samples follow a **t distribution** with $k-1$ **degrees of freedom** (here, $k=10$)
- Use **t-test** (or **Student's t-test**)
- **Null Hypothesis:** M_1 & M_2 are the same
- If we can **reject** null hypothesis, then
 - we conclude that the difference between M_1 & M_2 is **statistically significant**
 - Chose model with lower error rate

25

Estimating Confidence Intervals: Classifier Models M_1 vs. M_2

- Suppose we have 2 classifiers, M_1 and M_2 , which one is better?
- Use 10-fold cross-validation to obtain $\overline{err}(M_1)$ and $\overline{err}(M_2)$
- These mean error rates are just *estimates* of error on the true population of *future* data cases
- What if the difference between the 2 error rates is just attributed to *chance*?
 - Use a **test of statistical significance**
 - Obtain **confidence limits** for our error estimates

24

Estimating Confidence Intervals: t-test

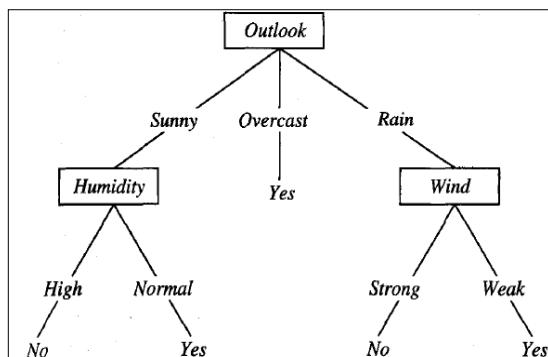
- If only 1 test set available: **pairwise comparison**
 - For i^{th} round of 10-fold cross-validation, the same cross partitioning is used to obtain $err(M_1)_i$ and $err(M_2)_i$
 - Average over 10 rounds to get $\overline{err}(M_1)$ and $\overline{err}(M_2)$
 - **t-test** computes **t-statistic** with $k-1$ **degrees of freedom**:
$$t = \frac{\overline{err}(M_1) - \overline{err}(M_2)}{\sqrt{\overline{var}(M_1 - M_2)/k}} \quad \text{where}$$
$$\overline{var}(M_1 - M_2) = \frac{1}{k} \sum_{i=1}^k [err(M_1)_i - err(M_2)_i - (\overline{err}(M_1) - \overline{err}(M_2))]^2$$
- If two test sets available: use **non-paired t-test**

$$\text{where } var(M_1 - M_2) = \sqrt{\frac{var(M_1)}{k_1} + \frac{var(M_2)}{k_2}},$$

where k_1 & k_2 are # of cross-validation samples used for M_1 & M_2 , resp.

26

Decision Tree learning: ID3



Decision trees represent a disjunction of conjunctions of constraints on the attribute values of instances.

$$\begin{aligned} & (\text{Outlook} = \text{Sunny} \wedge \text{Humidity} = \text{Normal}) \\ \vee & \quad (\text{Outlook} = \text{Overcast}) \\ \vee & \quad (\text{Outlook} = \text{Rain} \wedge \text{Wind} = \text{Weak}) \end{aligned}$$

For "Yes"

- Decision tree representation

- Objective

- Classify the instances by sorting them down the tree from the root to some leaf node, which provides the class of the instance
 - Each node tests the attribute of an instance whereas each branch descending from that node specifies one possible value of this attribute

- Testing

- Start at the root node, test the attribute specified by this node and identify the appropriate branch then repeat the process

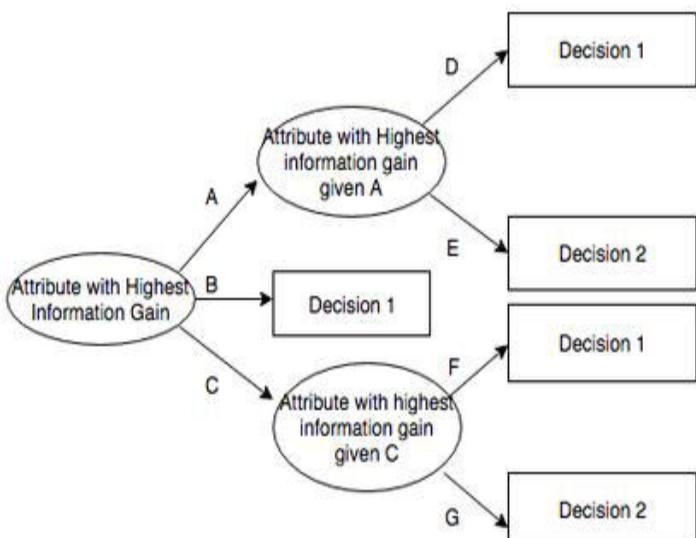
- Example

- Decision tree for classifying the Saturday mornings to test whether it is suitable for playing tennis based on the weather attributes

- The characteristics of best suited problems for DT

- Instances that are represented by attribute-value pairs
 - Each attribute takes small number of disjoint possible values (e.g., Hot, Mild, Cold) → discrete
 - Real-valued attributes like temperature → continuous
- The target function has discrete output values
 - Binary, multi-class and real-valued outputs
- Disjunctive descriptions may be required
- The training data may contain errors
 - DT robust to both errors: errors in classification and attributes
- The training data may contain missing attribute values

- Basics of DT learning - Iterative Dichotomiser 3 - ID3



Information gain

Information gain $IG(A)$ is the measure of the difference in entropy from before to after the set S is split on an attribute A . In other words, how much uncertainty in S was reduced after splitting set S on attribute A .

$$IG(A, S) = H(S) - \sum_{t \in T} p(t)H(t)$$

Where,

- $H(S)$ – Entropy of set S
- T – The subsets created from splitting set S by attribute A such that $S = \bigcup_{t \in T} t$
- $p(t)$ – The proportion of the number of elements in t to the number of elements in set S
- $H(t)$ – Entropy of subset t

In ID3, information gain can be calculated (instead of entropy) for each remaining attribute. The attribute with the **largest** information gain is used to split the set S on this iteration.

Entropy

Entropy $H(S)$ is a measure of the amount of uncertainty in the (data) set S (i.e. entropy characterizes the (data) set S).

$$H(S) = \sum_{x \in X} -p(x) \log_2 p(x)$$

Where,

- S – The current (data) set for which entropy is being calculated (changes every iteration of the ID3 algorithm)
- X – Set of classes in S
- $p(x)$ – The proportion of the number of elements in class x to the number of elements in set S

When $H(S) = 0$, the set S is perfectly classified (i.e. all elements in S are of the same class).

In ID3, entropy is calculated for each remaining attribute. The attribute with the **smallest** entropy is used to split the set S on this iteration. The higher the entropy, the higher the potential to improve the classification here.

- ID3 algorithm

- ID3 algorithm begins with the original set S as the root node.
- On each iteration of the algorithm, **it iterates through every unused attribute of the set S and calculates the entropy $H(S)$ (or information gain $IG(S)$) of that attribute.**
- It then selects the attribute which has the smallest entropy (or largest information gain) value.
 - The set S is then split by the selected attribute (e.g. age is less than 50, age is between 50 and 100, age is greater than 100) to produce subsets of the data.
- **The algorithm continues to recurs on each subset, considering only attributes never selected before.**

- Recursion on a subset may stop, When
 - All the elements in the class belong to same class
 - All instances does not belong to same class but there is no attribute to select
 - There is no example in the subset
- Steps in ID3
 - Calculate the entropy of every attribute using the data set S
 - Split the set S into subsets using the attribute for which the resulting entropy (after splitting) is minimum (or, equivalently, information gain is maximum)
 - Make a decision tree node containing that attribute
 - Recurs on subsets using remaining attributes.

ID3(Examples, Target_attribute, Attributes)

Examples are the training examples. Target_attribute is the attribute whose value is to be predicted by the tree. Attributes is a list of other attributes that may be tested by the learned decision tree. Returns a decision tree that correctly classifies the given Examples.

- Create a Root node for the tree
- If all Examples are positive, Return the single-node tree Root, with label = +
- If all Examples are negative, Return the single-node tree Root, with label = -
- If Attributes is empty, Return the single-node tree Root, with label = most common value of Target_attribute in Examples
- Otherwise Begin
 - A \leftarrow the attribute from Attributes that best* classifies Examples
 - The decision attribute for Root \leftarrow A
 - For each possible value, v_i , of A,
 - Add a new tree branch below Root, corresponding to the test $A = v_i$
 - Let Examples $_{v_i}$ be the subset of Examples that have value v_i for A
 - If Examples $_{v_i}$ is empty
 - Then below this new branch add a leaf node with label = most common value of Target_attribute in Examples
 - Else below this new branch add the subtree ID3(Examples $_{v_i}$, Target_attribute, Attributes - {A}))
- End
- Return Root

• Root - Significance

– Root node

- Which attribute should be tested at the root of the DT?
 - Decided based on the information gain or entropy
 - Significance
 - » The best attribute that classifies the training samples very well
 - » The attribute that should be tested in all the instances of the dataset
- @root node
 - Information gain is more or entropy has the least value
- Entropy measures the homogeneity or impurities in the instances
- Information gain measures the expected reduction in entropy

DT by example

Day	Outlook	Temperature	Humidity	Wind	PlayTennis	Target
D1	Sunny	Hot	High	Weak	No	
D2	Sunny	Hot	High	Strong	No	
D3	Overcast	Hot	High	Weak	Yes	
D4	Rain	Mild	High	Weak	Yes	
D5	Rain	Cool	Normal	Weak	Yes	
D6	Rain	Cool	Normal	Strong	No	
D7	Overcast	Cool	Normal	Strong	Yes	
D8	Sunny	Mild	High	Weak	No	
D9	Sunny	Cool	Normal	Weak	Yes	
D10	Rain	Mild	Normal	Weak	Yes	
D11	Sunny	Mild	Normal	Strong	Yes	
D12	Overcast	Mild	High	Strong	Yes	
D13	Overcast	Hot	Normal	Weak	Yes	
D14	Rain	Mild	High	Strong	No	

Entropy using frequency table of independent attributes on single dependent attribute
 (Target Variable here: PlayGolf)

Entropy using frequency table of single attribute on single attribute

$$E(S) = \sum_{i=1}^c -p_i \log_2 p_i$$

Play Golf	
Yes	No
9	5

$$\begin{aligned} \text{Entropy(PlayGolf)} &= \text{Entropy}(5,9) \\ &= \text{Entropy}(0.36, 0.64) \\ &= -(0.36 \log_2 0.36) - (0.64 \log_2 0.64) \\ &= 0.94 \end{aligned}$$

$$Gain(S, A) \equiv Entropy(S) - \sum_{v \in Values(A)} \frac{|S_v|}{|S|} Entropy(S_v)$$

Training examples: **9 yes / 5 no**

Day	Outlook	Humidity	Wind	Play
D1	Sunny	High	Weak	No
D2	Sunny	High	Strong	No
D3	Overcast	High	Weak	Yes
D4	Rain	High	Weak	Yes
D5	Rain	Normal	Weak	Yes
D6	Rain	Normal	Strong	No
D7	Overcast	Normal	Strong	Yes
D8	Sunny	High	Weak	No
D9	Sunny	Normal	Weak	Yes
D10	Rain	Normal	Weak	Yes
D11	Sunny	Normal	Strong	Yes
D12	Overcast	High	Strong	Yes
D13	Overcast	Normal	Weak	Yes
D14	Rain	High	Strong	No

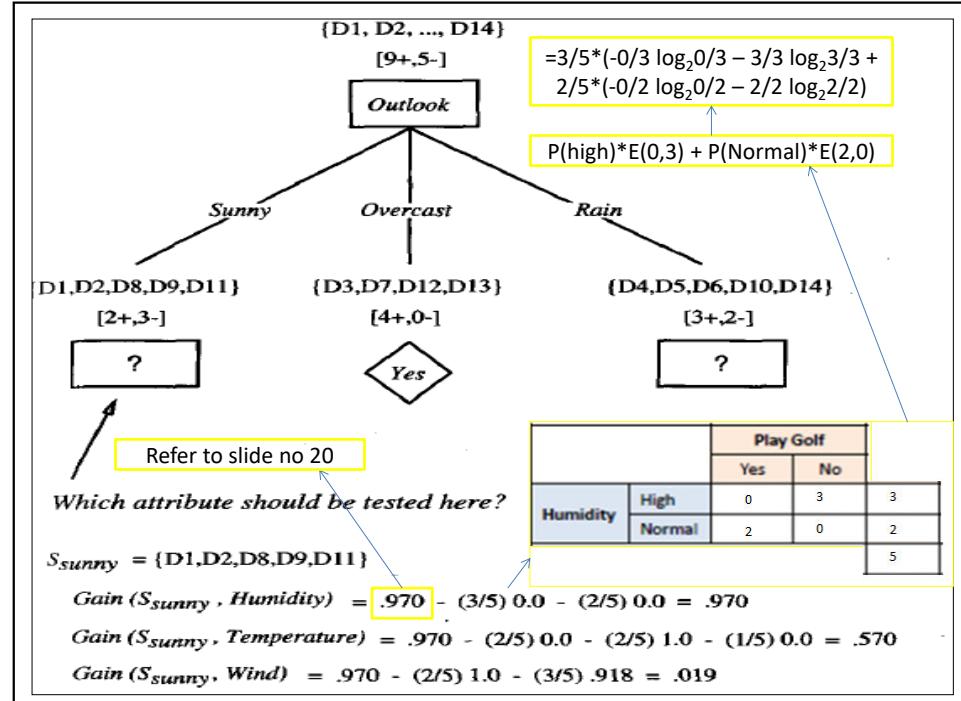
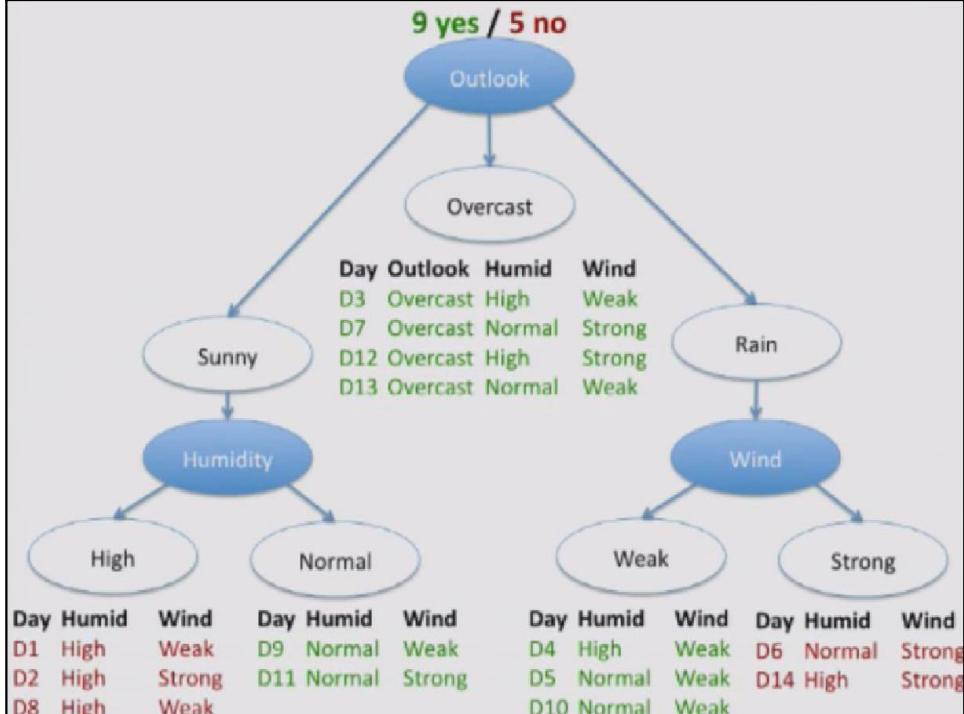
Entropy using frequency table of single attribute on two attribute

PlayGolf on Outlook
 PlayGolf on Temperature
 PlayGolf on Humidity
 PlayGolf on Windy

$$Gain(S, A) \equiv Entropy(S) - \sum_{v \in Values(A)} \frac{|S_v|}{|S|} Entropy(S_v)$$

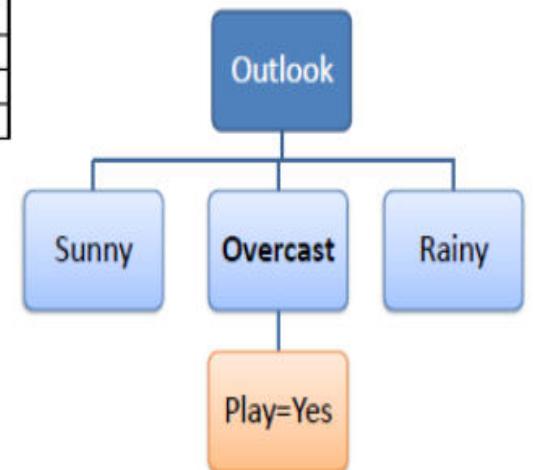
Choose the other nodes below the root node

Iterate the same procedure for finding the root node



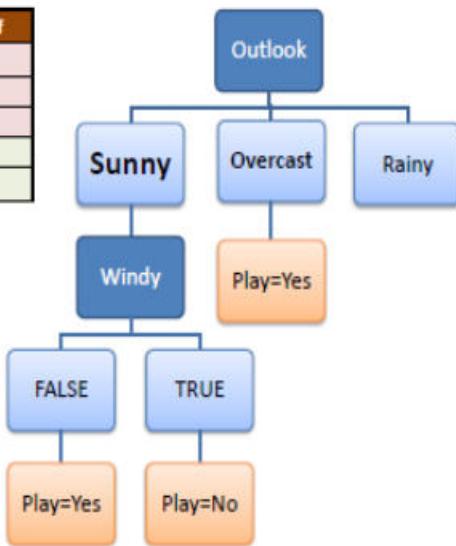
Branch with entropy of '0' is a leaf node

Temp	Humidity	Windy	Play Golf
Hot	High	FALSE	Yes
Cool	Normal	TRUE	Yes
Mild	High	TRUE	Yes
Hot	Normal	FALSE	Yes



Branch with entropy more than '0' needs further splitting

Temp	Humidity	Windy	Play Golf
Mild	High	FALSE	Yes
Cool	Normal	FALSE	Yes
Mild	Normal	FALSE	Yes
Cool	Normal	TRUE	No
Mild	High	TRUE	No



Try this

Age	Income	Student	Credit_rating	Buys_computer
<=30	high	no	fair	no
<=30	high	no	excellent	no
31...40	high	no	fair	yes
>40	medium	no	fair	yes
>40	low	yes	fair	yes
>40	low	yes	excellent	no
31...40	low	yes	excellent	yes
<=30	medium	no	fair	no
<=30	low	yes	fair	yes
>40	medium	yes	fair	yes
<=30	medium	yes	excellent	yes
31...40	medium	no	excellent	yes
31...40	high	yes	fair	yes
>40	medium	no	excellent	no

R₁: IF (Outlook=Sunny) AND (Windy=FALSE) THEN Play=Yes

R₂: IF (Outlook=Sunny) AND (Windy=TRUE) THEN Play=No

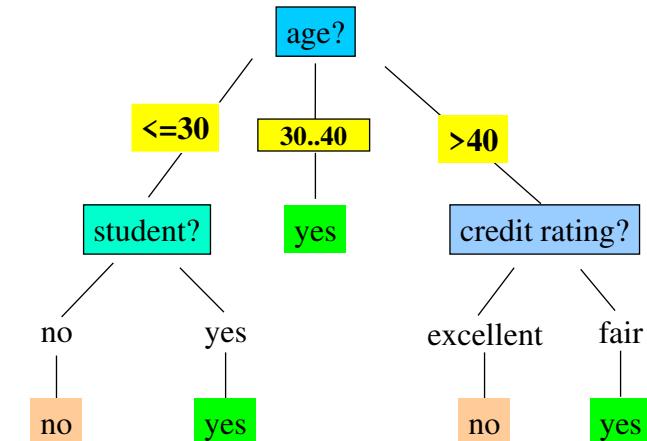
R₃: IF (Outlook=Overcast) THEN Play=Yes

R₄: IF (Outlook=Rainy) AND (Humidity=High) THEN Play=No

R₅: IF (Outlook=Rainy) AND (Humidity=Normal) THEN Play=Yes



Output: A Decision Tree for "buys_computer"



- **HYPOTHESIS SPACE SEARCH IN DECISION TREE LEARNING**
- ID3 can be characterized as searching a space of hypotheses for one that fits the training examples
 - The hypothesis space searched by ID3 is the set of possible decision trees
 - simple-to complex, hill-climbing search through this hypothesis space,
 - beginning with the empty tree, then considering progressively more elaborate hypotheses in search of a decision tree that correctly classifies the training data
- ID3 capabilities and limitations by considering search space and search strategy
 - ID3 hypothesis space of all decision trees is a complete space of finite discrete-valued functions, relative to the available attributes
 - Adv – avoids search incomplete hypothesis space that might not contain target function
 - ID3 maintains only a single current hypothesis as it searches through the space of decision trees
 - In candidate-elimination algorithms set of all hypotheses consistent with the available training examples
 - it does not have the ability to determine how many alternative decision trees are consistent with the available training data
 - ID3 does not support backtracking leads to local optimum.
 - ID3 uses all training examples at each step in the search to make statistically based decisions regarding how to refine its current hypothesis

- **INDUCTIVE BIAS IN DECISION TREE LEARNING**
 - inductive bias is the set of assumptions that, together with the training data, deductively justify the classifications assigned by the learner to future instances
 - selects in favor of shorter trees over longer ones, and
 - selects trees that place the attributes with highest information gain closest to the root.
 - A closer approximation to the inductive bias of ID3: Shorter trees are preferred over longer trees. Trees that place high information gain attributes close to the root are preferred over those that do not.
- Inductive bias ID3 versus CANDIDATE-ELIMINATION
 - ID3 → Preference bias
 - searches a **complete hypothesis space**
 - searches **incompletely** through this space, **from simple to complex hypotheses**, until it finds a hypothesis consistent with the data and it does not go beyond that.
 - **Inductive bias:** a consequence of the ordering of hypotheses by its search strategy and its hypothesis space introduces no additional bias
 - » inductive bias of ID3 follows from its search strategy **thus a preference for certain hypotheses over others** (e.g., for shorter hypotheses)
 - CANDIDATE-ELIMINATION → Restriction bias
 - searches **an incomplete hypothesis space** (i.e., one that can **express only a subset of the potentially teachable concepts**)
 - searches this space **completely**, **finding every hypothesis consistent** with the training data
 - **Inductive bias:** consequence of the expressive power of its hypothesis representation. Its search strategy introduces no additional bias.
 - » Solely relies on search space representation

• Why Prefer Short Hypotheses?

– Occam's razor

- Prefer the simplest hypothesis that fits the data

– Claim

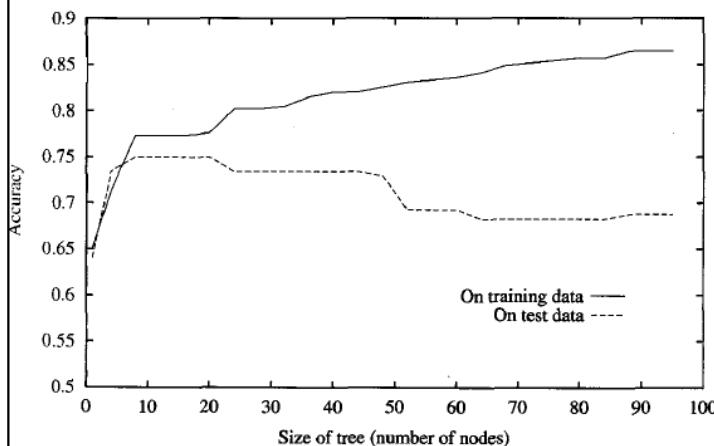
- short hypothesis that coincidentally fits the training data

– Critics

- Not always true
- Internal representation strategy of learners decides the size of the hypothesis.
 - Two learners arrive at different hypotheses, both justifying their contradictory conclusions by Occam's razor
- Occam's razor will produce two different hypotheses from the same training examples

Criterion is to be used to determine the correct final tree size are:

- Assess the utility of post-pruning nodes using separate set of examples different from training
- Assess the improvement by applying expansion or pruning
- Apply explicit condition for tree growth



• ISSUES IN DECISION TREE LEARNING

– Avoid Overfitting of data

- **Definition:** Given a hypothesis space H , a hypothesis $h \in H$ is said to overfit the training data if there exists some alternative hypothesis $h' \in H$, such that h has smaller error than h' over the training examples, but h' has a smaller error than h over the entire distribution of instances.

• Reasons for overfitting

- when there is noise in the data, or when the number of training examples is too small to produce a representative sample of the true target function

• Methods to avoid overfitting

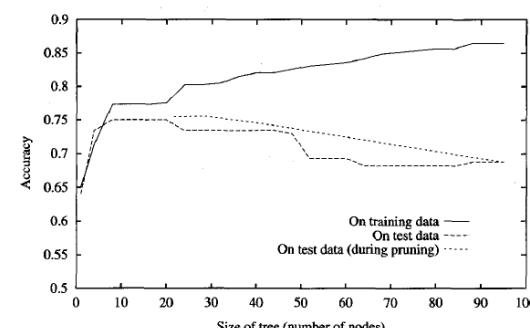
- Stop growing the tree earlier
- Firstly allow to overfit and then post-prune the tree (more successful)

• REDUCED ERROR PRUNING

– How exactly might we use a validation set to prevent overfitting?

- Pruning a decision node consists of removing the subtree rooted at that node, making it a leaf node, and assigning it the most common classification of the training examples affiliated with that node.

- Condition: Nodes are removed only if the resulting pruned tree performs no worse than the original over the validation set



- RULE POST-PRUNING

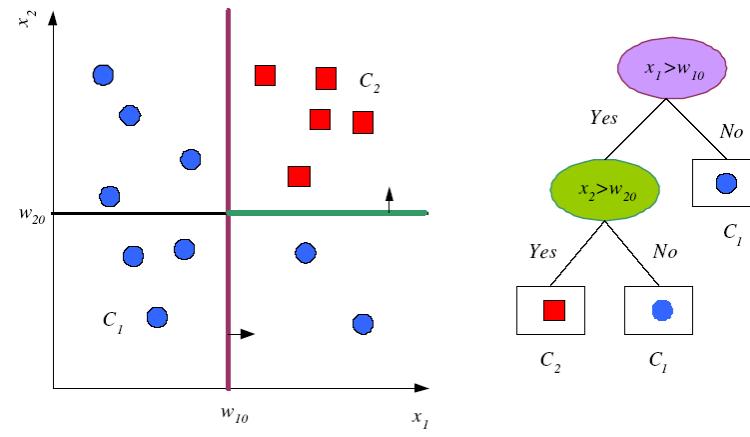
1. Infer the decision tree from the training set, growing the tree until the training data is fit as well as possible and allowing overfitting to occur.
2. Convert the learned tree into an equivalent set of rules by creating one rule for each path from the root node to a leaf node.
3. Prune (generalize) each rule by removing any preconditions that result in improving its estimated accuracy.
4. Sort the pruned rules by their estimated accuracy, and consider them in this sequence when classifying subsequent instances.

Divide and Conquer

- Internal decision nodes
 - Univariate: Uses a single attribute, x_i
 - Numeric x_i : Binary split : $x_i > w_m$
 - Discrete x_i : n -way split for n possible values
 - Multivariate: Uses all attributes, \mathbf{x}
- Leaves
 - Classification: Class labels, or proportions
 - Regression: Numeric; r average, or local fit
- Learning is greedy; find the best split recursively (Breiman et al, 1984; Quinlan, 1986, 1993)

43

Tree Uses Nodes, and Leaves



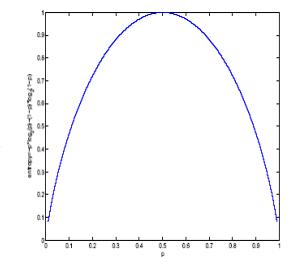
42

Classification Trees (ID3, CART, C4.5)

- For node m , N_m instances reach m , N_m^i belong to C_i

$$\hat{P}(C_i | \mathbf{x}, m) \equiv p_m^i = \frac{N_m^i}{N_m}$$
- Node m is pure if p_m^i is 0 or 1
- Measure of impurity is entropy

$$I_m = -\sum_{i=1}^K p_m^i \log_2 p_m^i$$



44

Best Split

- If node m is pure, generate a leaf and stop, otherwise split and continue recursively
- Impurity after split: N_{mj} of N_m take branch j . N^i_{mj} belong to C_i

$$\hat{P}(C_i | \mathbf{x}, m, j) \equiv p_{mj}^i = \frac{N_{mj}^i}{N_{mj}}$$

$$I_m = -\sum_{j=1}^n \frac{N_{mj}}{N_m} \sum_{i=1}^K p_{mj}^i \log_2 p_{mj}^i$$

- Find the variable and split that min impurity (among all variables -- and split positions for numeric variables)

45

- **Classification:** split the dataset based on homogeneity of data
- **Regression: (There is no class label)**
 - » fit a regression model to the target variable using each of the independent variables
 - » for each independent variable, the **data is split at several split points**
 - » At each split point, the "error" between the predicted value and the actual values is squared to get a "Sum of Squared Errors (SSE)".
 - » split point errors across the variables are compared and the variable /point yielding the lowest SSE is chosen as the root node/split point.
 - » This process is recursively continued

GenerateTree(\mathcal{X})

```

If NodeEntropy( $\mathcal{X}$ ) <  $\theta_I$  /* eq. 9.3
  Create leaf labelled by majority class in  $\mathcal{X}$ 
  Return
i  $\leftarrow$  SplitAttribute( $\mathcal{X}$ )
For each branch of  $\mathbf{x}_i$ 
  Find  $\mathcal{X}_i$  falling in branch
  GenerateTree( $\mathcal{X}_i$ )
SplitAttribute( $\mathcal{X}$ )
MinEnt  $\leftarrow$  MAX
For all attributes  $i = 1, \dots, d$ 
  If  $\mathbf{x}_i$  is discrete with  $n$  values
    Split  $\mathcal{X}$  into  $\mathcal{X}_1, \dots, \mathcal{X}_n$  by  $\mathbf{x}_i$ 
    e  $\leftarrow$  SplitEntropy( $\mathcal{X}_1, \dots, \mathcal{X}_n$ ) /* eq. 9.8 */
    If e < MinEnt MinEnt  $\leftarrow$  e; bestf  $\leftarrow$  i
  Else /*  $\mathbf{x}_i$  is numeric */
    For all possible splits
      Split  $\mathcal{X}$  into  $\mathcal{X}_1, \mathcal{X}_2$  on  $\mathbf{x}_i$ 
      e  $\leftarrow$  SplitEntropy( $\mathcal{X}_1, \mathcal{X}_2$ )
    If e < MinEnt MinEnt  $\leftarrow$  e; bestf  $\leftarrow$  i
Return bestf

```

Regression Trees

- Error at node m :

$$b_m(\mathbf{x}) = \begin{cases} 1 & \text{if } \mathbf{x} \in \mathcal{X}_m : \mathbf{x} \text{ reaches node } m \\ 0 & \text{otherwise} \end{cases}$$

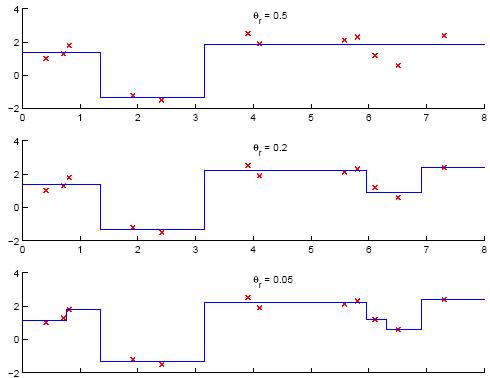
$$E_m = \frac{1}{N_m} \sum_t (r^t - g_m)^2 b_m(\mathbf{x}^t) \quad g_m = \frac{\sum_t b_m(\mathbf{x}^t) r^t}{\sum_t b_m(\mathbf{x}^t)}$$

- After splitting:

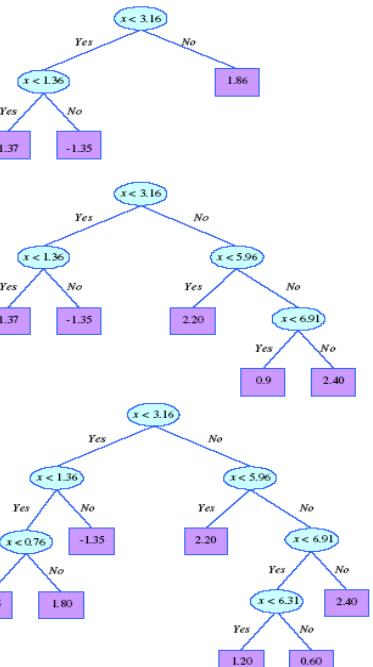
$$b_{mj}(\mathbf{x}) = \begin{cases} 1 & \text{if } \mathbf{x} \in \mathcal{X}_{mj} : \mathbf{x} \text{ reaches node } m \text{ and branch } j \\ 0 & \text{otherwise} \end{cases}$$

$$E_m = \frac{1}{N_m} \sum_j \sum_t (r^t - g_{mj})^2 b_{mj}(\mathbf{x}^t) \quad g_{mj} = \frac{\sum_t b_{mj}(\mathbf{x}^t) r^t}{\sum_t b_{mj}(\mathbf{x}^t)}$$

Model Selection in Trees:

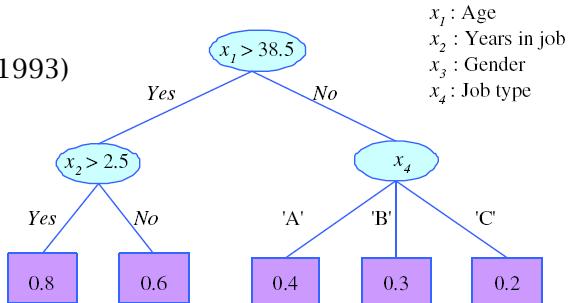


49



Rule Extraction from Trees

C4.5Rules
(Quinlan, 1993)



- R1: IF (age>38.5) AND (years-in-job>2.5) THEN $y = 0.8$
- R2: IF (age>38.5) AND (years-in-job≤2.5) THEN $y = 0.6$
- R3: IF (age≤38.5) AND (job-type='A') THEN $y = 0.4$
- R4: IF (age≤38.5) AND (job-type='B') THEN $y = 0.3$
- R5: IF (age≤38.5) AND (job-type='C') THEN $y = 0.2$

51

Pruning Trees

- Remove subtrees for better generalization (decrease variance)
 - Prepruning: Early stopping
 - Postpruning: Grow the whole tree then prune subtrees which overfit on the pruning set
- Prepruning is faster, postpruning is more accurate (requires a separate pruning set)

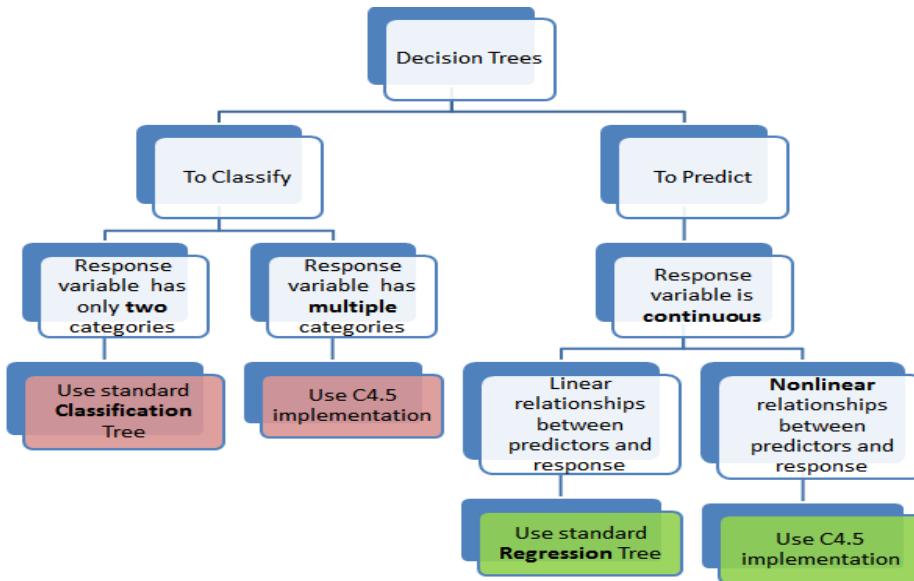
50

Learning Rules

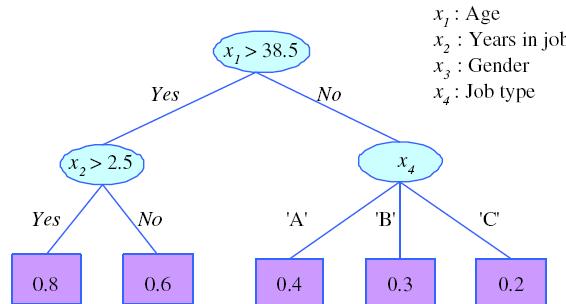
- Rule induction is similar to tree induction but
 - tree induction is breadth-first,
 - rule induction is depth-first; one rule at a time
- Rule set contains rules; rules are conjunctions of terms
- Rule covers an example if all terms of the rule evaluate to true for the example
- Sequential covering: Generate rules one at a time until all positive examples are covered
- IREP (Fürnkranz and Widmer, 1994), Ripper (Cohen, 1995)

52

Classification and Regression Trees (CART)



Rule Extraction from Trees



- R1: IF (age>38.5) AND (years-in-job>2.5) THEN $y = 0.8$
- R2: IF (age>38.5) AND (years-in-job≤2.5) THEN $y = 0.6$
- R3: IF (age≤38.5) AND (job-type='A') THEN $y = 0.4$
- R4: IF (age≤38.5) AND (job-type='B') THEN $y = 0.3$
- R5: IF (age≤38.5) AND (job-type='C') THEN $y = 0.2$

55

Comparison between ID3, C4.5 and CART

	Splitting Criteria	Attribute type	Missing values	Pruning Strategy	Outlier Detection
ID3	Information Gain	Handles only Categorical value	Do not handle missing values.	No pruning is done	Susceptible to outliers
CART	Towing Criteria	Handles both Categorical & Numeric value	Handle missing values.	Cost-Complexity pruning is used	Can handle Outliers
C4.5	Gain Ratio	Handles both Categorical & Numeric value	Handle missing values.	Error Based pruning is used	Susceptible to outliers

Regression Tree

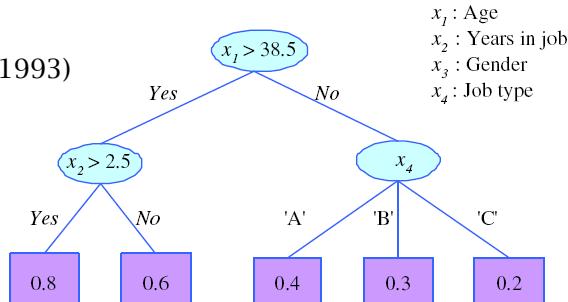
Predictors

Outlook	Temp.	Humidity	Windy	Hours Played
Rainy	Hot	High	False	26
Rainy	Hot	High	True	30
Overcast	Hot	High	False	46
Sunny	Mild	High	False	46
Sunny	Cool	Normal	False	62
Sunny	Cool	Normal	True	23
Overcast	Cool	Normal	True	43
Rainy	Mild	High	False	36
Rainy	Cool	Normal	False	38
Sunny	Mild	Normal	False	46
Rainy	Mild	Normal	True	48
Overcast	Mild	High	True	62
Overcast	Hot	Normal	False	44
Sunny	Mild	High	True	30

Decision Tree learning: CART

Rule Extraction from Trees

C4.5Rules
(Quinlan, 1993)



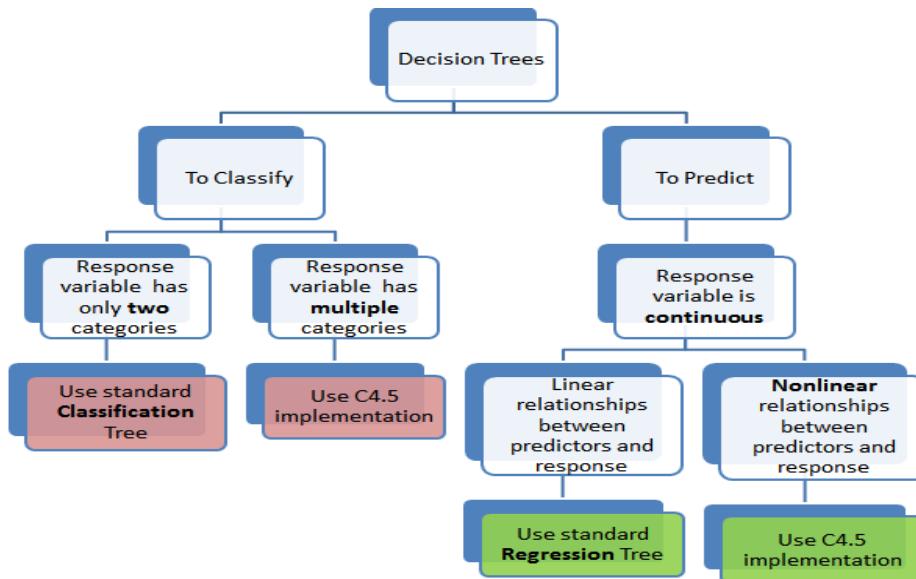
- R1: IF (age>38.5) AND (years-in-job>2.5) THEN $y = 0.8$
- R2: IF (age>38.5) AND (years-in-job≤2.5) THEN $y = 0.6$
- R3: IF (age≤38.5) AND (job-type='A') THEN $y = 0.4$
- R4: IF (age≤38.5) AND (job-type='B') THEN $y = 0.3$
- R5: IF (age≤38.5) AND (job-type='C') THEN $y = 0.2$

- **Classification:** split the dataset based on homogeneity of data
- **Regression:** (There is no class label)
 - » fit a regression model to the target variable using each of the independent variables
 - » for each independent variable, the **data is split at several split points**
 - » At each split point, the "error" between the predicted value and the actual values is squared to get a "Sum of Squared Errors (SSE)".
 - » split point errors across the variables are compared and the variable /point yielding the lowest SSE is chosen as the root node/split point.
 - » This process is recursively continued

Learning Rules

- Rule induction is similar to tree induction but
 - tree induction is breadth-first,
 - rule induction is depth-first; one rule at a time
- Rule set contains rules; rules are conjunctions of terms
- Rule covers an example if all terms of the rule evaluate to true for the example
- Sequential covering: Generate rules one at a time until all positive examples are covered
- IREP (Fürnkranz and Widmer, 1994), Ripper (Cohen, 1995)

Classification and Regression Trees (CART)

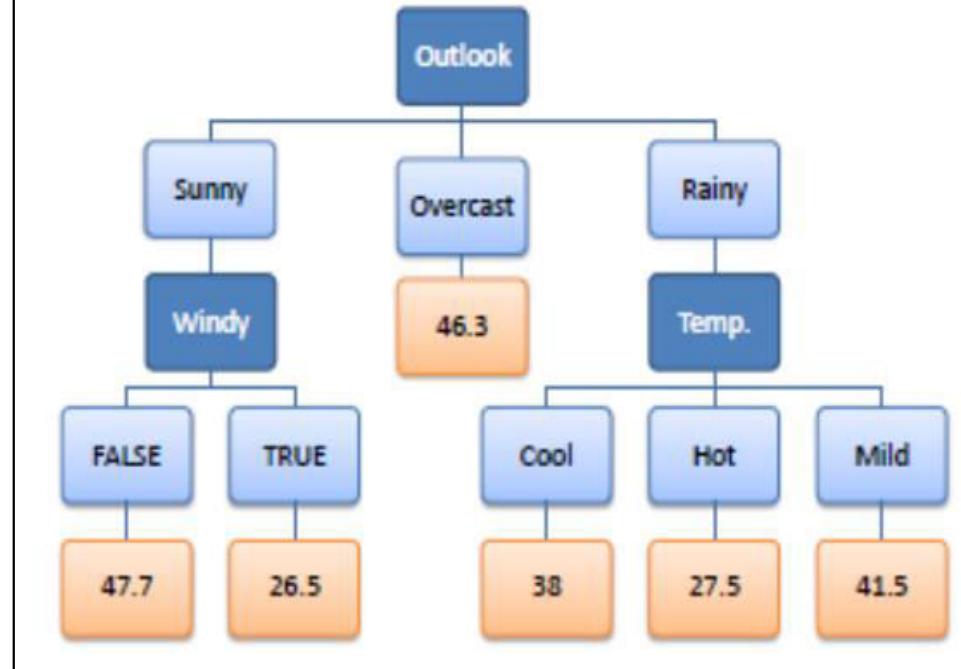


Regression Tree

Predictors				Target
Outlook	Temp.	Humidity	Windy	Hours Played
Rainy	Hot	High	False	26
Rainy	Hot	High	True	30
Overcast	Hot	High	False	48
Sunny	Mild	High	False	46
Sunny	Cool	Normal	False	62
Sunny	Cool	Normal	True	23
Overcast	Cool	Normal	True	43
Rainy	Mild	High	False	36
Rainy	Cool	Normal	False	38
Sunny	Mild	Normal	False	46
Rainy	Mild	Normal	True	48
Overcast	Mild	High	True	62
Overcast	Hot	Normal	False	44
Sunny	Mild	High	True	30

Comparison between ID3, C4.5 and CART

	Splitting Criteria	Attribute type	Missing values	Pruning Strategy	Outlier Detection
ID3	Information Gain	Handles only Categorical value	Do not handle missing values.	No pruning is done	Susceptible to outliers
CART	Gini Criteria	Handles both Categorical & Numeric value	Handle missing values.	Cost-Complexity pruning is used	Can handle Outliers
C4.5	Gain Ratio	Handles both Categorical & Numeric value	Handle missing values.	Error Based pruning is used	Susceptible to outliers



Standard deviation

Golf players = {25, 30, 46, 45, 52, 23, 43, 35, 38, 46, 48, 52, 44, 30}

Average of golf players = $(25 + 30 + 46 + 45 + 52 + 23 + 43 + 35 + 38 + 46 + 48 + 52 + 44 + 30) / 14 = 39.78$

Standard deviation of golf players = $\sqrt{[(25 - 39.78)^2 + (30 - 39.78)^2 + (46 - 39.78)^2 + \dots + (30 - 39.78)^2] / 14} = 9.32$

Outlook

Outlook can be sunny, overcast and rain. We need to calculate standard deviation of golf players for all of these outlook candidates.

Sunny outlook

Day	Outlook	Temp.	Humidity	Wind	Golf Players
1	Sunny	Hot	High	Weak	25
2	Sunny	Hot	High	Strong	30
8	Sunny	Mild	High	Weak	35
9	Sunny	Cool	Normal	Weak	38
11	Sunny	Mild	Normal	Strong	48

Golf players for sunny outlook = {25, 30, 35, 38, 48}

Average of golf players for sunny outlook = $(25+30+35+38+48)/5 = 35.2$

Standard deviation of golf players for sunny outlook = $\sqrt{[(25 - 35.2)^2 + (30 - 35.2)^2 + \dots + (48 - 35.2)^2] / 5} = 7.78$

Overcast outlook

Day	Outlook	Temp.	Humidity	Wind	Golf Players
3	Overcast	Hot	High	Weak	46
7	Overcast	Cool	Normal	Strong	43
12	Overcast	Mild	High	Strong	52
13	Overcast	Hot	Normal	Weak	44

Golf players for overcast outlook = {46, 43, 52, 44}

Average of golf players for overcast outlook = $(46 + 43 + 52 + 44) / 4 = 46.25$

Standard deviation of golf players for overcast outlook = $\sqrt{[(46 - 46.25)^2 + (43 - 46.25)^2 + \dots + (52 - 46.25)^2] / 4} = 3.49$

Rainy outlook

Day	Outlook	Temp.	Humidity	Wind	Golf Players
4	Rain	Mild	High	Weak	45
5	Rain	Cool	Normal	Weak	52
6	Rain	Cool	Normal	Strong	23
10	Rain	Mild	Normal	Weak	46
14	Rain	Mild	High	Strong	30

Golf players for rainy outlook = {45, 52, 23, 46, 30}

Average of golf players for rainy outlook = $(45 + 52 + 23 + 46 + 30) / 5 = 39.2$

Standard deviation of golf players for rainy outlook = $\sqrt{[(45 - 39.2)^2 + (52 - 39.2)^2 + \dots + (30 - 39.2)^2] / 5} = 10.87$

Summarizing standard deviations for the outlook feature

Outlook	Stdev of Golf Players	Instances
Overcast	3.49	4
Rain	10.87	5
Sunny	7.78	5

Weighted standard deviation for outlook = $(4/14) \times 3.49 + (5/14) \times 10.87 + (5/14) \times 7.78 = 7.66$

You might remember that we have calculated the global standard deviation of golf players 9.32 in previous steps. Standard deviation reduction is difference of the global standard deviation and standard deviation for current feature. In this way, maximized standard deviation reduction will be the decision node.

Standard deviation reduction for outlook = $9.32 - 7.66 = 1.66$

Temperature

Temperature can be hot, cool or mild. We will calculate standard deviations for those candidates.

Hot temperature

Day	Outlook	Temp.	Humidity	Wind	Golf Players
1	Sunny	Hot	High	Weak	25
2	Sunny	Hot	High	Strong	30
3	Overcast	Hot	High	Weak	46
13	Overcast	Hot	Normal	Weak	44

Golf players for hot temperature = {25, 30, 46, 44}

Standard deviation of golf players for hot temperature = 8.95

Cool temperature

Day	Outlook	Temp.	Humidity	Wind	Golf Players
5	Rain	Cool	Normal	Weak	52
6	Rain	Cool	Normal	Strong	23
7	Overcast	Cool	Normal	Strong	43
9	Sunny	Cool	Normal	Weak	38

Golf players for cool temperature = {52, 23, 43, 38}

Standard deviation of golf players for cool temperature = 10.51

Mild temperature

Day	Outlook	Temp.	Humidity	Wind	Golf Players
4	Rain	Mild	High	Weak	45
8	Sunny	Mild	High	Weak	35
10	Rain	Mild	Normal	Weak	46
11	Sunny	Mild	Normal	Strong	48
12	Overcast	Mild	High	Strong	52
14	Rain	Mild	High	Strong	30

Golf players for mild temperature = {45, 35, 46, 48, 52, 30}

Standard deviation of golf players for mild temperature = 7.65

Wind

Wind is a binary class, too. It can either be Strong or Weak.

Strong Wind

Day	Outlook	Temp.	Humidity	Wind	Golf Players
2	Sunny	Hot	High	Strong	30
6	Rain	Cool	Normal	Strong	23
7	Overcast	Cool	Normal	Strong	43
11	Sunny	Mild	Normal	Strong	48
12	Overcast	Mild	High	Strong	52
14	Rain	Mild	High	Strong	30

Golf players for strong wind= {30, 23, 43, 48, 52, 30}

Standard deviation for golf players for strong wind = 10.59

Weak Wind

1	Sunny	Hot	High	Weak	25
3	Overcast	Hot	High	Weak	46
4	Rain	Mild	High	Weak	45
5	Rain	Cool	Normal	Weak	52
8	Sunny	Mild	High	Weak	35
9	Sunny	Cool	Normal	Weak	38
10	Rain	Mild	Normal	Weak	46
13	Overcast	Hot	Normal	Weak	44

Golf players for weakk wind= {25, 46, 45, 52, 35, 38, 46, 44}

Standard deviation for golf players for weak wind = 7.87

Summarizing standard deviations for wind feature

Wind	Stdev of Golf Player	Instances
Strong	10.59	6
Weak	7.87	8

Weighted standard deviation for wind = $(6/14) \times 10.59 + (8/14) \times 7.87 = 9.03$

Standard deviation reduction for wind = $9.32 - 9.03 = 0.29$

So, we've calculated standard deviation reduction values for all features. The winner is outlook because it has the highest score.

Feature	Standard Deviation Reduction
Outlook	1.66
Temperature	0.47
Humidity	0.27
Wind	0.29

We'll put outlook decision at the top of decision tree. Let's monitor the new sub data sets for the candidate branches of outlook feature.

Sunny outlook and normal humidity

Day	Outlook	Temp.	Humidity	Wind	Golf Players
9	Sunny	Cool	Normal	Weak	38
11	Sunny	Mild	Normal	Strong	48

Standard deviation for sunny outlook and normal humidity = 5

Summarizing standard deviations for humidity feature when outlook is sunny

Humidity	Stdev for Golf Players	Instances
High	4.08	3
Normal	5.00	2

Weighted standard deviations for sunny outlook and humidity = $(3/5) \times 4.08 + (2/5) \times 5 = 4.45$

Standard deviation reduction for sunny outlook and humidity = $7.78 - 4.45 = 3.33$

Sunny outlook and Strong Wind

Day	Outlook	Temp.	Humidity	Wind	Golf Players
2	Sunny	Hot	High	Strong	30
11	Sunny	Mild	Normal	Strong	48

Standard deviation for sunny outlook and strong wind = 9

Sunny outlook and Weak Wind

Day	Outlook	Temp.	Humidity	Wind	Golf Players
1	Sunny	Hot	High	Weak	25
8	Sunny	Mild	High	Weak	35
9	Sunny	Cool	Normal	Weak	38

Standard deviation for sunny outlook and weak wind = 5.56

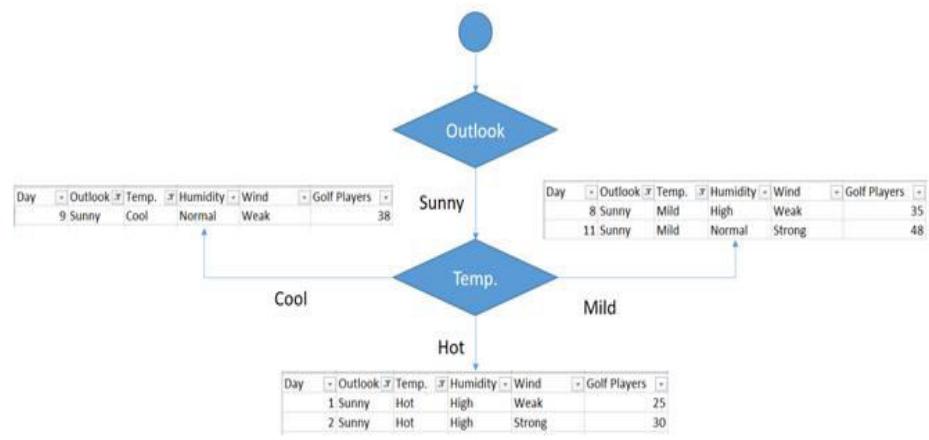
Wind	Stdev for Golf Players	Instances
Strong	9	2
Weak	5.56	3

Weighted standard deviations for sunny outlook and wind = $(2/5) \times 9 + (3/5) \times 5.56 = 6.93$

Standard deviation reduction for sunny outlook and wind = $7.78 - 6.93 = 0.85$

We've calculated standard deviation reductions for sunny outlook. The winner is temperature.

Feature	Standard Deviation Reduction
Temperature	4.18
Humidity	3.33
Wind	0.85



Putting temperature decision at the bottom of sunny outlook

Pruning

Cool branch has one instance in its sub data set. We can say that if outlook is sunny and temperature is cool, then there would be 38 golf players. But what about hot branch? There are still 2 instances. Should we add another branch for weak wind and strong wind? No, we should not. Because this causes over-fitting. We should terminate building branches, for example if there are less than five instances in the sub data set. Or standard deviation of the sub data set can be less than 5% of the entire data set. I prefer to apply the first one. I will terminate the branch if there are less than 5 instances in the current sub data set. If this termination condition is satisfied, then I will calculate the average of the sub data set. This operation is called as pruning in decision tree trees.

Overcast outlook

Overcast outlook branch has already 4 instances in the sub data set. We can terminate building branches for this leaf. Final decision will be average of the following table for overcast outlook.

Day	Outlook	Temp.	Humidity	Wind	Golf Players
3	Overcast	Hot	High	Weak	46
7	Overcast	Cool	Normal	Strong	43
12	Overcast	Mild	High	Strong	52
13	Overcast	Hot	Normal	Weak	44

If outlook is overcast, then there would be $(46+43+52+44)/4 = 46.25$ golf players.

Rainy Outlook

Day	Outlook	Temp.	Humidity	Wind	Golf Players
4	Rain	Mild	High	Weak	45
5	Rain	Cool	Normal	Weak	52
6	Rain	Cool	Normal	Strong	23
10	Rain	Mild	Normal	Weak	46
14	Rain	Mild	High	Strong	30

We need to find standard deviation reduction values for the rest of the features in same way for the sub data set above.

Standard deviation for rainy outlook = 10.87

Notice that we will use this value as global standard deviation for this branch in reduction step.

Rainy outlook and temperature

Temperature	Standard deviation for golf players	instances
Cool	14.50	2
Mild	7.32	3

Weighted standard deviation for rainy outlook and temperature = $(2/5) \times 14.50 + (3/5) \times 7.32 = 10.19$

Standard deviation reduction for rainy outlook and temperature = $10.87 - 10.19 = 0.67$

Rainy outlook and humidity

Humidity	Standard deviation for golf players	instances
High	7.50	2
Normal	12.50	3

Weighted standard deviation for rainy outlook and humidity = $(2/5) \times 7.50 + (3/5) \times 12.50 = 10.50$

Standard deviation reduction for rainy outlook and humidity = $10.87 - 10.50 = 0.37$

Rainy outlook and wind

Wind	Standard deviation for golf players	instances
Weak	3.09	3
Strong	3.5	2

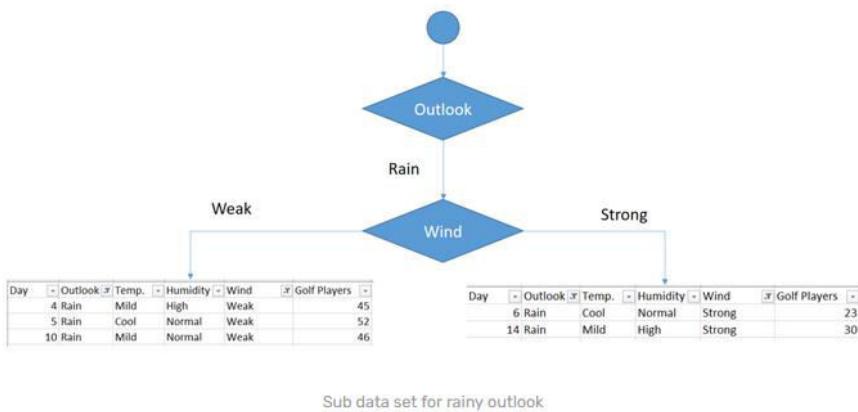
Weighted standard deviation for rainy outlook and wind = $(3/5) \times 3.09 + (2/5) \times 3.5 = 3.25$

Standard deviation reduction for rainy outlook and wind = $10.87 - 3.25 = 7.62$

Summarizing rainy outlook

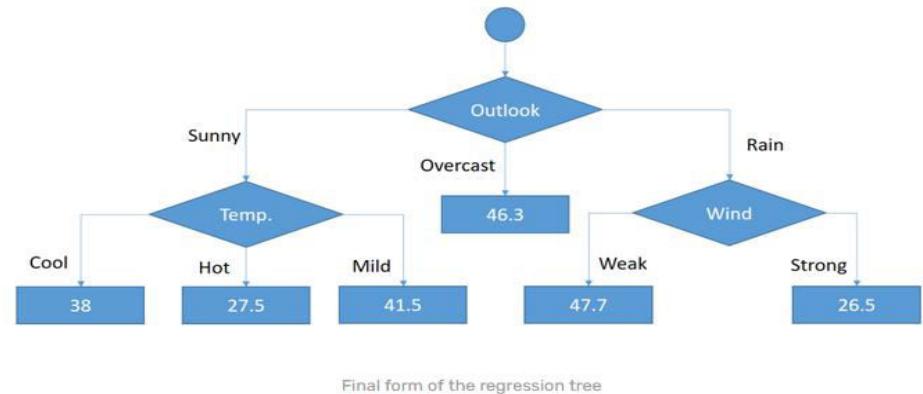
As illustrated below, the winner is wind feature.

Feature	Standard deviation reduction
Temperature	0.67
Humidity	0.37
Wind	7.62



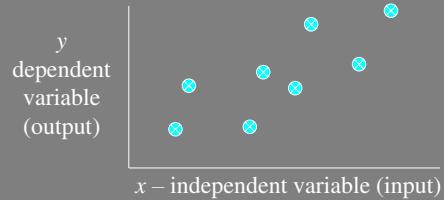
As seen, both branches have items less than 5. Now, we can terminate these leafs based on the termination rule.

So, Final form of the decision tree is demonstrated below.



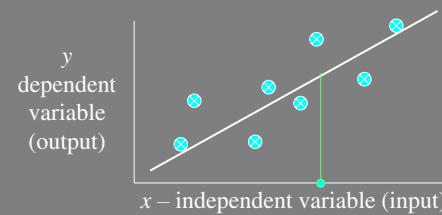
Regression

- For classification the output(s) is nominal
- In regression the output is continuous
 - Function Approximation
- Many models could be used – Simplest is linear regression
 - Fit data with the best hyper-plane which "goes through" the points



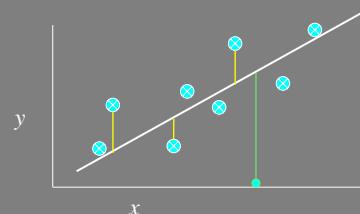
Regression

- For classification the output(s) is nominal
- In regression the output is continuous
 - Function Approximation
- Many models could be used – Simplest is linear regression
 - Fit data with the best hyper-plane which "goes through" the points



Regression

- For classification the output(s) is nominal
- In regression the output is continuous
 - Function Approximation
- Many models could be used – Simplest is linear regression
 - Fit data with the best hyper-plane which "goes through" the points
 - For each point the difference between the predicted point and the actual observation is the *residue*



3

# CORRECT(x)	ATTITUDE (y)
17	94
13	73
12	59
15	80
16	93
14	85
16	66
16	79
18	77
19	91

Simple Linear Regression

- For now, assume just one (input) independent variable x , and one (output) dependent variable y
 - Multiple linear regression assumes an input vector \mathbf{x}
 - Multivariate linear regression assumes an output vector \mathbf{y}
- We "fit" the points with a line (i.e. hyper-plane)
- Which line should we use?
 - Choose an objective function
 - For simple linear regression we choose sum squared error (SSE)
 - $\sum (\text{predicted}_i - \text{actual}_i)^2 = \sum (\text{residue}_i)^2$
 - Thus, find the line which minimizes the sum of the squared residues (e.g. least squares)
 - This mimics the case where data points were sampled from the actual hyperplane with Gaussian noise added

4

Linear Regression function

$$y = a + bx$$

slope (b) of Regression Line

$$b = r \frac{s_y}{s_x}$$

y -Intercept (a) of Regression Line

$$a = \bar{y} - b \bar{x}$$

Pearson Correlation Coefficient (r)

$$r = \frac{\sum ((x - \bar{x})(y - \bar{y}))}{\sqrt{\sum (x - \bar{x})^2 \sum (y - \bar{y})^2}}$$

# CORRECT(x)	ATTITUDE(y)
17	94
13	73
12	59
15	80
16	93
14	85
16	66
16	79
18	77
19	91

#CORRECT(x)	ATTITUDE(y)	(x - \bar{x})	(y - \bar{y})	(x - \bar{x})(y - \bar{y})	(x - \bar{x}) ²	(y - \bar{y}) ²
17	94	1.4	14.3	20.02	1.96	204.49
13	73	-2.6	-6.7	17.42	6.76	44.89
12	59	-3.6	-20.7	74.52	12.96	428.49
15	80	-0.6	0.3	-0.18	0.36	0.09
16	93	0.4	13.3	5.32	0.16	176.89
14	85	-1.6	5.3	-8.48	2.56	28.09
16	66	0.4	-13.7	-5.48	0.16	187.69
16	79	0.4	-0.7	-0.28	0.16	0.49
18	77	2.4	-2.7	-6.48	5.76	7.29
19	91	3.4	11.3	38.42	11.56	127.69
$\bar{x} = 15.6$		$\bar{y} = 79.7$		$\Sigma = 134.8$	$\Sigma = 42.4$	$\Sigma = 1206.1$

$$r = \frac{134.8}{\sqrt{42.4 \times 1206.1}} = 0.596$$

$$S_y = \sqrt{\frac{\sum (y - \bar{y})^2}{n-1}}$$

$$= \sqrt{\frac{1206.1}{9}} = 11.576$$

$$S_x = \sqrt{\frac{\sum (x - \bar{x})^2}{n-1}}$$

$$= \sqrt{\frac{42.4}{9}} = 2.171$$

SSE and Linear Regression

- SSE chooses to square the difference of the predicted vs actual. Why square?
- Don't want residues to cancel each other
- Could use absolute or other distances to solve problem
 - $\Sigma |predicted_i - actual_i|$: L1 vs L2
- SSE leads to a parabolic error surface which is good for gradient descent
- Which line would least squares choose?
 - There is always one “best” fit



15

Sum-squared Error (SSE)

$$SSE = \sum y_{observed} - y_{predicted})^2$$

$$TSS = \sum y_{observed} - \bar{y}_{observed})^2$$

$$R^2 = 1 - \frac{SSE}{TSS}$$

What is Best Fit?

- The smaller the SSE, the better the fit
- Hence,
 - Linear regression attempts to minimize SSE (or similarly to maximize R²)
- Assume 2 dimensions

$$Y = \beta_0 + \beta_1 X$$

Analytical Solution

$$\beta_0 = \frac{\bar{y} - \beta_1 \bar{x}}{n}$$

$$\beta_1 = \frac{n \bar{xy} - \bar{x} \bar{y}}{n \bar{x^2} - (\bar{x})^2}$$

Example (I)

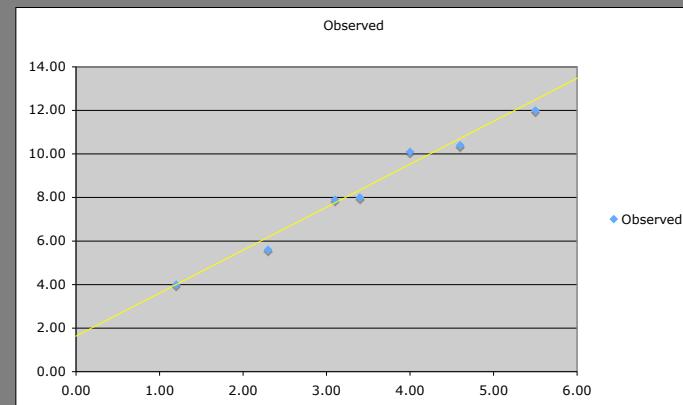
x	y	x^2	xy
1.20	4.00	1.44	4.80
2.30	5.60	5.29	12.88
3.10	7.90	9.61	24.49
3.40	8.00	11.56	27.20
4.00	10.10	16.00	40.40
4.60	10.40	21.16	47.84
5.50	12.00	30.25	66.00
24.10	58.00	95.31	223.61

$$\begin{aligned} b_1 &= \frac{n \sum xy - \bar{x} \bar{y}}{n \sum x^2 - (\bar{x})^2} \\ &= \frac{7 \cdot 223.61 - 24.10 \cdot 58.00}{7 \cdot 95.31 - 24.10^2} \\ &= \frac{1565.27 - 1397.80}{667.17 - 580.81} \\ &= \frac{167.47}{86.36} = \underline{\underline{1.94}} \end{aligned}$$

$$\begin{aligned} b_0 &= \bar{y} - b_1 \bar{x} \\ &= \frac{58.00 - 1.94 \cdot 24.10}{7} \\ &= \frac{11.27}{7} = \underline{\underline{1.61}} \end{aligned}$$

Target: $y=2x+1.5$

Example (II)



Example (III)

x	y (obs)	y (pred)	SSE	TSS
1.20	4.00	3.94	0.004	18.367
2.30	5.60	6.07	0.221	7.213
3.10	7.90	7.62	0.078	0.149
3.40	8.00	8.21	0.044	0.082
4.00	10.10	9.37	0.533	3.292
4.60	10.40	10.53	0.017	4.470
5.50	12.00	12.28	0.078	13.796
			0.975	47.369

$$R^2 = 1 - \frac{SSE}{TSS} = 1 - \frac{0.975}{47.369} = 0.98$$

let's say that we had data on the prices of homes on sale and the actual number of sales of homes:

Price(thousands of \$) Sales of new homes

x	y
160	126
180	103
200	82
220	75
240	82
260	40
280	20

Plot the data and check for the linear relationship between attributes if any? and find the least square regression line and compute SSE.