Table of Contents

		Dauction	
	1.1.	What is the World Wide Web?	·· 1
	1.2.	A Brief History of the Web and the Internet	2
	1.3.	Web Data Mining	·· 4 ·· 6 ·· 6
	1.4.		
	1.5.	How to Read this Book ······	
	Bibli	ographic Notes ·····	12
Pa	art I:	Data Mining Foundations	
2.	Ass	ociation Rules and Sequential Patterns	
		ociation Rules and Sequential Patterns ************************************	13
	2.1.	Basic Concepts of Association Rules	13 16 16
	2.1. 2.2.	Basic Concepts of Association Rules Apriori Algorithm 2.2.1. Frequent Itemset Generation 2.2.2 Association Rule Generation	13 16 16 20
	2.1.2.2.2.3.	Apriori Algorithm	13 16 16 20 22 22 24 26

	2.6.	Basic (Concepts of Sequential Patterns	· 37
	2.7.	Mining 2.7.1. 2.7.2.	Sequential Patterns Based on GSP GSP Algorithm Mining with Multiple Minimum Supports	. 39
	2.8.	Mining 2.8.1. 2.8.2.	Sequential Patterns Based on PrefixSpan ······ PrefixSpan Algorithm ······· Mining with Multiple Minimum Supports ······	· 46
	2.9.	2.9.1. 2.9.2.	ating Rules from Sequential Patterns Sequential Rules Label Sequential Rules Class Sequential Rules	· 50
	Biblio	ographi	c Notes ·····	· 52
3.	Sup	ervise	d Learning ······	• 55
	3.1.	Basic (Concepts	. 55
	3.2.	3.2.1. 3.2.2.	Handling of Continuous Attributes	· 62 · 63 · 67
	3.3.	Classif 3.3.1. 3.3.2.	ier Evaluation Evaluation Methods Precision, Recall, F-score and Breakeven Point	· 71
	3.4.	Rule Ir 3.4.1. 3.4.2. 3.4.3.	nduction	· 75 · 78
	3.5.	Classif 3.5.1. 3.5.2. 3.5.3.	ication Based on Associations	· 82 · 86
	3.6.	Naïve	Bayesian Classification ······	· 87
	3.7.	Naïve 3.7.1. 3.7.2. 3.7.3.	Bayesian Text Classification Probabilistic Framework Naïve Bayesian Model Discussion	· 92 · 93
	3.8.		rt Vector Machines ······ Linear SVM: Separable Case ······	

			Table of Contents	VI
		3.8.2.	Linear SVM: Non-Separable Case ·····	
		3.8.3.	Nonlinear SVM: Kernel Functions	
			rest Neighbor Learning ·····	
	3.10.		nble of Classifiers ·····	
			Bagging Boosting	
	Biblio		ic Notes ·····	
4.	Uns	uperv	ised Learning ······	• 117
	4.1.	Basic	Concepts ·····	· 117
	4.2.	K-mea	ans Clustering ·····	120
			K-means Algorithm	
		4.2.2.	5	
		4.2.3.	Strengths and Weaknesses ·····	· 124
	4.3.		sentation of Clusters	
			Common Ways of Representing Clusters	
		4.3.2	Clusters of Arbitrary Shapes	
	4.4.		chical Clustering	
		4.4.1.		
			Complete-Link Method ······ Average-Link Method ······	
		4.4.4.	Strengths and Weaknesses ······	
	4.5		ice Functions	
		4.5.1.		
		4.5.2.	Binary and Nominal Attributes	
		4.5.3.	Text Documents ·····	
	4.6.	Data S	Standardization	· 139
	4.7.	Handl	ing of Mixed Attributes	· 141
	4.8.	Which	Clustering Algorithm to Use?	· 143
	4.9.	Cluste	r Evaluation ·····	· 143
	4.10.	Discov	vering Holes and Data Regions ·····	· 146
	Biblio	ographi	ic Notes ·····	· 149
5.	Part	ially S	upervised Learning	151
	5.1.		ing from Labeled and Unlabeled Examples	
		511	FM Algorithm with Naïve Bayesian Classification	· 153

		5.1.2. 5.1.3. 5.1.4. 5.1.5. 5.1.6.	Co-Training Self-Training Transductive Support Vector Machines Graph-Based Methods Discussion	158 159
	5.2.	Learni 5.2.1. 5.2.2. 5.2.3. 5.2.4. 5.2.5.	ng from Positive and Unlabeled Examples Applications of PU Learning Theoretical Foundation Building Classifiers: Two-Step Approach Building Classifiers: Direct Approach Discussion	165 168 169 175
	Appe	endix: D	Perivation of EM for Naïve Bayesian Classification ···	179
			c Notes ·····	
Pa	rt II	: W (eb Mining	
6.	Info	rmatio	n Retrieval and Web Search	183
6.			n Retrieval and Web Search	
6.	Info 6.1. 6.2.	Basic	Concepts of Information Retrieval ation Retrieval Models Boolean Model Vector Space Model Statistical Language Model	184 187 188 188
6.	6.1.	Basic (Inform 6.2.1. 6.2.2. 6.2.3.	Concepts of Information Retrieval ation Retrieval Models Boolean Model Vector Space Model	184 187 188 188 191
6.	6.1. 6.2.	Basic (Inform 6.2.1. 6.2.2. 6.2.3. Releva	Concepts of Information Retrieval ation Retrieval Models Boolean Model Vector Space Model Statistical Language Model ance Feedback ation Measures	184 187 188 188 191 192 195
6.	6.1.6.2.6.3.	Basic (Inform 6.2.1. 6.2.2. 6.2.3. Releva	Concepts of Information Retrieval ation Retrieval Models Boolean Model Vector Space Model Statistical Language Model ance Feedback	184 187 188 181 191 192 195 199 200 200

	6.7.	Latent	Semantic Indexing	215
		6.7.1.		215
		6.7.2.		
		6.7.3.		
		6.7.4.		
	6.8.	Web S	earch ·····	222
	6.9.	Meta-S	Search: Combining Multiple Rankings	225
		6.9.1.	Combination Using Similarity Scores ······	
		6.9.2.	Combination Using Rank Positions	227
	6.10.		pamming	
			Content Spamming ·····	
		6.10.2.	Link Spamming	231
			Hiding Techniques	
			Combating Spam ·····	
	Bibli	ographi	c Notes ·····	235
7	1 :1-	A a ls .	a:a	007
7.	LINK	Anaiy	sis	237
	7.1.	Social	Network Analysis	238
		7.1.1	Centrality ·····	
		7.1.2	Prestige	241
	7.2.	Co-Cita	ation and Bibliographic Coupling	243
		7.2.1.	Co-Citation ·····	244
		7.2.2.	Bibliographic Coupling	
	7.3.	PageR	ank ·····	
		7.3.1.	PageRank Algorithm ·····	
		7.3.2.	Strengths and Weaknesses of PageRank	253
		7.3.3.	•	
	7.4.			
			HITS Algorithm ·····	
		7.4.2.		259
		7.4.3.	Relationships with Co-Citation and Bibliographic	050
		711	Coupling	
		7.4.4.	Strengths and Weaknesses of HITS	
	7.5.		unity Discovery	
		7.5.1. 7.5.2.	Problem Definition	
		7.5.2. 7.5.3.	Maximum Flow Communities	
		7.5.3. 7.5.4.	Email Communities Based on Betweenness	
		7.5. 4 . 7.5.5.	Overlapping Communities of Named Entities	
			C.C. app g Commando of Hamea Endico	0

	Biblio	ographic	Notes ·····	271		
8.	Web	Web Crawling ·····				
	8.1.	A Basi 8.1.1. 8.1.2.	c Crawler Algorithm ····································	275		
	8.2.	8.2.1. 8.2.2. 8.2.3. 8.2.4. 8.2.5.	Petching Parsing Stopword Removal and Stemming Link Extraction and Canonicalization Spider Traps Page Repository Concurrency	277 278 280 280 282 283		
	8.3.	Univer 8.3.1. 8.3.2.	sal CrawlersScalability Coverage vs Freshness vs Importance	286		
	8.4.	Focuse	ed Crawlers ·····	289		
	8.5.	Topica 8.5.1. 8.5.2. 8.5.3.	I Crawlers	294 300		
	8.6.	Evalua	tion ·····	310		
	8.7.	Crawle	er Ethics and Conflicts ·····	315		
	8.8.	Some	New Developments	318		
	Biblio	ographic	Notes ·····	320		
9.	Stru	ctured	Data Extraction: Wrapper Generation ·	323		
	9.1	Prelimi 9.1.1. 9.1.2. 9.1.3.	inaries Two Types of Data Rich Pages Data Model HTML Mark-Up Encoding of Data Instances	324 326		
	9.2.	Wrapp 9.2.1. 9.2.2. 9.2.3. 9.2.4.	er Induction ·····	330 330 333 337		

	9.3.	Instance-Based Wrapper Learning	338
	9.4.	Automatic Wrapper Generation: Problems	342
	9.5.	String Matching and Tree Matching	344
	9.6.	Multiple Alignment	350
	9.7.	Building DOM Trees	356
	9.8.	Extraction Based on a Single List Page: Flat Data Records	
		9.8.1. Two Observations about Data Records	358
		9.8.3. Identifying Data Records in Data Regions	
		9.8.4. Data Item Alignment and Extraction	365
		9.8.5. Making Use of Visual Information	366
		9.8.6. Some Other Techniques	366
	9.9.	Extraction Based on a Single List Page: Nested Data Records	367
	9.10.	Extraction Based on Multiple Pages	
		9.10.1. Using Techniques in Previous Sections	373
	0.44	9.10.2. RoadRunner Algorithm	
	9.11.	Some Other Issues	
		9.11.2. Disjunction or Optional ······	376
		9.11.3. A Set Type or a Tuple Type	377
		9.11.4. Labeling and Integration	378
	0.40	9.11.5. Domain Specific Extraction	
		Discussion ····	
	Biblic	ographic Notes ·····	379
10.	Info	mation Integration ·····	381
	10.1.	Introduction to Schema Matching	382
	10.2.	Pre-Processing for Schema Matching	384
	10.3.	Schema-Level Match ·····	385

	10.3.1. Linguistic Approaches	385 386
10.4.	Domain and Instance-Level Matching	387
10.5.		
10.6.	1:m Match	391
10.7.	Some Other Issues ·····	392
	10.7.1. Reuse of Previous Match Results ······	
	10.7.2. Matching a Large Number of Schemas ······	393
10.8		
10.0.	10.8.1. A Clustering Based Approach	397
	10.8.2. A Correlation Based Approach	400
10.9.		406
	10.9.1. Structural Appropriateness and the	406
	10.9.2. Lexical Appropriateness ······	408
	10.9.3. Instance Appropriateness·····	409
Biblio	graphic Notes ·····	410
Opini	ion Mining	411
11.1.		
	11.1.3. Classification Using a Score Function	
11.2.	Feature-Based Opinion Mining and Summarization ··	417
		424
	of Format 1	425
	11.2.4. Feature Extraction from Reviews of	
11.3.	Comparative Sentence and Relation Mining	432
	11.2.1 Droblem Definition	422
	11.3.1. Problem Definition	433
	10.5. 10.6. 10.7. 10.8. 10.9. Bibliog	10.3.2. Constraint Based Approaches 10.4. Domain and Instance-Level Matching 10.5. Combining Similarities 10.6. 1:m Match 10.7. Some Other Issues 10.7.1. Reuse of Previous Match Results 10.7.2. Matching a Large Number of Schemas 10.7.3. Schema Match Results 10.7.4. User Interactions 10.8. Integration of Web Query Interfaces 10.8.1. A Clustering Based Approach 10.8.2. A Correlation Based Approach 10.8.3. An Instance Based Approach 10.9. Constructing a Unified Global Query Interface 10.9.1. Structural Appropriateness and the Merge Algorithm 10.9.2. Lexical Appropriateness 10.9.3. Instance Appropriateness 10.9.3. Instance Appropriateness 10.9.1. Classification Using Text Classification Methods 11.1.1. Classification Using Text Classification Methods 11.1.2. Classification Using a Score Function 11.2.1. Problem Definition 11.2.2. Object Feature Extraction 11.2.3. Feature Extraction from Pros and Cons of Format 1 11.2.4. Feature Extraction from Reviews of of Formats 2 and 3 11.2.5. Opinion Orientation Classification Mining 11.3. Comparative Sentence and Relation Mining

			Table of Contents	XIII
	11.5.	Opinio Opinio 11.5.1. 11.5.2. 11.5.3. 11.5.4.	Extraction of Comparative Relations n Search n Spam Objectives and Actions of Opinion Spamming Types of Spam and Spammers Hiding Techniques Spam Detection Notes	439 441 441 442 443 444
12.			Mining	
	12.1.	12.1.1	Sollection and Pre-Processing Sources and Types of Data Key Elements of Web Usage Data Pre-Processing	452
	12.2	Data M	lodeling for Web Usage Mining	462
	12.3	12.3.1. 12.3.2. 12.3.3 12.3.4	ery and Analysis of Web Usage Patterns Session and Visitor Analysis Cluster Analysis and Visitor Segmentation Association and Correlation Analysis Analysis of Sequential and Navigational Patterns Classification and Prediction Based on Web User Transactions	466 467 471 475
	12.4.	Discus	sion and Outlook ······	482
	Biblio	graphic	Notes ·····	482
Ref	erenc	es ·····		485
Inde	ex			517