Contents

	Preface Contributors 2 Part I KNOWLEDGE DISCOVERY AND DATA MINING IN THEORY				
Pa					
1	Estimating concept difficulty with cross entropy.				
	K. N	azar and M.A. Bramer	3		
	1.1	Introduction	3		
		1.1.1 Attributes and examples	4		
	1.2	The need for bias in machine learning	4		
		1.2.1 Bias interference	5		
	1.3	Concept dispersion and feature interaction	7		
		Why do we need to estimate concept difficulty?	8		
	1.5	Data-based difficulty measures	10		
		1.5.1 μ-ness	10		
		1.5.2 Variation	11		
		1.5.3 Blurring	11		
	1.6	Why use the J measure as a basis for a blurring measure?	14		
	1.7	Effects of various data characteristics on blurring and variation	15		
		1.7.1 Irrelevant attributes	17		
	1.8	Dataset analysis	19		
	1.9	Problems with information-theoretic-based blurring measures	21		
	1.10	Estimating attribute relevance with the RELIEF algorithm	22		
		1.10.1 Experiments with RELIEFF	23		
	1.11	Blurring over disjoint subsets	26		
		Results and discussion	27		
	1.13	Summary and future research	28		
	1.14	References	30		
2	Analysing outliers by searching for plausible hypotheses.				
		iu and G. Cheng	32		
	2.1	Introduction	32		
	2.2	Statistical treatment of outliers	33		

vi Contents

2.3 An algorithm for outlier analysis	
2.4 Experimental results	35
2.4.1 Case I: onchocerciasis	38
2.4.2 Case II: glaucoma	38
2.5 Evaluation	39
2.6 Concluding remarks	41
2.7 Acknowledgments	44
2.8 References	44
	44
3 Attribute-value distribution as a technique for increasing the	
of data mining. D. McSherry	46
3.1 Introduction	40 46
3.2 Targeting a restricted class of rules	40 48
3.3 Discovery effort and yield	50
3.4 Attribute-value distribution	50 53
3.5 Experimental results	
3.5.1 The contact-lens data	5 5
3.5.2 The project-outcome dataset	55 60
3.0 Discussion and conclusions	60 62
3.7 References	63
	03
4 Using background knowledge with attribute-oriented data mining. M. Shabcatt, S. McClan, and B. Santa	
- Stapesti, S. Wicken and B. Scotney	64
4.1 Introduction	64
4.2 Partial value model	66
4.2.1 Definition: partial value	66
4.2.2 Definition: partial-value relation	66
4.2.3 Example	66
4.2.4 Aggregation	67
4.2.5 Definition: simple count operator	67
4.2.0 Example	68
4.2.7 Definition: simple aggregate operator	68
4.2.0 Example	68
4.2.9 Aggregation of partial values	69
4.2.10 Definition: partial value aggregate operator	69
1.2.11 Example	69
4.2.12 Definition: partial-value count operator	70
1.2.13 Theoretical framework	70
and database - the role of background	
Knowledge	71
	71
integrity constraints	73
simple companison predicate	73
4.3.4 Examples of simple comparison predicates	74
4.3.5 Definition: table-based predicate	74

			Con	ntents	vii		
		4.3.6	Example of a table-based predicate		74		
			Expression of rules as table-based predicates		74		
			Reengineering the database		74		
	4.4		attribute count operators		76		
			Example		77		
		4.4.2	Example		77		
			Example		78		
		4.4.4	Quasi-independence		79		
		4.4.5	Example		81		
		4.4.6	Interestingness		81		
		4.4.7	Example		82		
	4.5	Relate	ed work		82		
	4.6	Concl	lusions		83		
	4.7	Ackno	owledgments		84		
	4.8	Refer	ences		84		
5	A development framework for temporal data mining.						
		X. Chen and I. Petrounias					
			duction		87		
	5.2	Analy	sis and representation of temporal features		89		
			Time domain		89		
		5.2.2	Calendar expression of time		90		
			Periodicity of time		92		
		5.2.4	Time dimensions in temporal databases		94		
	5.3	Poten	tial knowledge and temporal data mining problems		95		
		5.3.1	Forms of potential temporal knowledge		95		
		5.3.2	Associating knowledge with temporal features		97		
			Temporal mining problems		98		
	5.4	A fran	nework for temporal data mining		100		
			A temporal mining language		100		
		5.4.2	System architecture		102		
	5.5	An ex	ample: discovery of temporal association rules		104		
			Mining problem		104		
		5.5.2	Description of mining tasks in TQML		106		
			Search algorithms		107		
	5.6	Concl	lusion and future research direction		110		
	5.7	Refere	ences		110		
6	An i	ntegrat	ed architecture for OLAP and data mining. Z. Cher	ı	114		
	6.1		duction		115		
	6.2	2 Preliminaries					
			Decision-support queries		116		
			Data warehousing		116		
			Basics of OLAP		117		
			Star schema		118		

	^	6.2.5 A materialised view for sales profit	118
	0	.3 Differences between OLAP and data mining	119
		6.3.1 Basic concepts of data mining	119
		6.3.2 Different types of query can be answered at different levels	
		6.3.3 Aggregation semantics	120
		6.3.4 Sensitivity analysis	120
		6.3.5 Different assumptions or heuristics may be needed at different levels	122
	6.	4 Combining OLAP and data mining: the feedback sandwich model	123
		6.4.1 Two different ways of combining OLAP and data mining	123
		6.4.2 The feedback sandwich model	124
	6.5	Towards integrated architecture	125
		Towards integrated architecture for combined OLAP/data mining	
	6.6		126
		6.6.1 On the use and war st	128
		6.6.1 On the use and reuse of intensional historical data	128
		6.6.2 How data mining can benefit OLAP 6.6.3 OLAP-enriched data mining	131
	6.7	Conclusion	133
		References	135
			135
F	art II	KNOWLEDGE DISCOVERY AND DATA MINING IN	
		PRACTICE PRACTICE	
			137
7	Em	pirical studies of the knowledge discovery approach to	
		M. Lloyd-Williams	
	,.1	introduction	139
	7.2	Knowledge discovery and data mining	139
		7.2.1 The knowledge discovery process	140
		7.2.2 Artificial neural networks	140
	7.3	Empirical studies	147
		7.3.1 The 'Health for all' database	149
		7.3.2 The 'Babies at risk of intrapartum and the state of	149
	. .	and unity databases	153
	7.4	Conclusions	155
	7.5	References	157
_			158
8	Dire	ct knowledge discovery and interpretation from a multilayer	
		on of Carrie, S.J. Idylor, IVI.A. Fow and A I D Form	100
		and oddiedolf	160
	8.2	The MLP network	160
	8.3	The low-back-pain MLP network	161
			163

	8.4	The i	nterpretation and knowledge-discovery method	163			
		8.4.1		163			
		8.4.2		164			
		8.4.3	, ,	164			
		8.4.4	Knowledge learned by the MLP from the training data	166			
		8.4.5	MLP network validation and verification	167			
	8.5	Know	vledge discovery from LBP example training cases	168			
		8.5.1	Discovery of the feature detectors for example				
			training cases	168			
		8.5.2	Discovery of the significant inputs for example				
			training cases	168			
		8.5.3	<u> </u>				
			training cases	168			
		8.5.4		170			
		8.5.5	0 1	170			
	8.6		eledge discovery from all LBP MLP training cases	173			
			Discussion of the class key input rankings	173			
	8.7		ation of the LBP LMP network	176			
			Validation of the training cases	176			
		8.7.2	Validation of the test cases	176			
	8.8		lusions	177			
			re work	177			
			owledgments	178			
	8.11	Refer	ences	178			
9	Discovering knowledge from low-quality meteorological databases.						
			rd and V.J. Rayward-Smith	180			
	9.1	Intro	duction	180			
		9.1.1	The meteorological domain	180			
	9.2	The p	preprocessing stage	182			
			Visualisation	182			
			Missing values	182			
			Unreliable data	185			
		9.2.4	Discretisation	186			
			Feature selection	187			
			Feature construction	188			
	9.3		lata-mining stage	188			
			Simulated annealing	189			
			SA with missing and unreliable data	190 195			
	9.4	12 toolkit for knowledge discovery					
	9.5		ts and analysis	196			
			Results from simulated annealing	198			
			Results from C5.0	198			
		9.5.3		199			
	9.6	Summ	nary	199			

x Contents

9.1	7 Discussion and further work		
9.8	References	200	
	References	201	
10 A	meteorological knowledge discovery		
A.	10 A meteorological knowledge-discovery environment. A.G. Büchner, J.C.L. Chan, S.L.Hung and J.G. Hughes		
10	10.1 Introduction		
10.	10.2 Some meteorological background		
	10.2.1 Available data sources	205	
	10.2.2 Related work	206	
10.	3 MADAME's architecture	208	
10.		211	
	10.4.1 The design	212	
	10.4.2 Information extraction	212	
	10.4.3 Data cleansing	213	
	10.4.4 Data processing	214	
	10.4.5 Data loading and refreshing	215	
10.5	The knowledge-discovery components	216	
	10.5.1 Knowledge modelling	217	
	10.5.2 Domain knowledge	217	
10.6	Prediction trial runs	221	
	10.6.1 Nowcasting of heavy rainfall	222	
	10.6.2 Landslide nowcasting	222	
10.7	Conclusions and further work	223	
10.8	Acknowledgments	224	
10.9	References	224	
		225	
11 Mining the organic compound jungle – a functional programming			
11.1	Introduction	227	
11.2	Decision-support requirements in the pharmaceutical	227	
	industry industry		
	11.2.1 Graphical comparison	227	
	11.2.2 Structural keys	228	
	11.2.3 Fingerprints	228	
	11.2.4 Variable-sized fingerprints	229	
11.3	Functional programming language Gofer	230	
	11.3.1 Functional programming	230	
11.4	Design of prototype tool and main functions	231	
	11.4.1 Design of tool	234	
	11.4.2 Main functions	234	
11.5	Methodology	235	
11.6	Results	236	
	11.6.1 Sets A, B and C (256, 512 and 1024 bytes)	236	
	11.0.2 Set D (2048 bytes)	237	
11.7	Conclusions	237	
		238	

			Contents	хi
	11.8	Future work		239
		References		239
	11.5	References		433
12	Data :	mining with neural networks – an applied example in		
	unde	rstanding electricity consumption patterns.		
	P. Bri	P. Brierley and B. Batty		
	12.1	Neural networks		241
		12.1.1 What are neural networks?		241
		12.1.2 Why use neural networks?		242
		12.1.3 How do neural networks process information?		242
		12.1.4 Things to be aware of		244
	12.2	Electric load modelling		250
		12.2.1 The data being mined		250
		12.2.2 Why forecast electricity demand?		252
	,	12.2.3 Previous work		253
		12.2.4 Network used		254
	12.3	Total daily load model		255
		12.3.1 Overfitting and generalisation		264
	12.4	Rule extraction		266
		12.4.1 Day of the week		267
		12.4.2 Time of year		268
		12.4.3 Growth		269
		12.4.4 Weather factors		271
		12.4.5 Holidays		272
	12.5	Model comparisons		274
	12.6	Half-hourly model		276
	•	12.6.1 Initial input data		276
		12.6.2 Results		277
		12.6.3 Past loads		286
		12.6.4 How the model is working		286
		12.6.5 Extracting the growth		287
		12.6.6 Populations of models		288
	12.7	Summary		288
	12.8	References		289
	12.9 Appendixes			292
		12.9.1 Backpropagation weight update rule		292
		12.9.2 Fortran 90 code for a multilayer perceptron		299
Ind	lex			304