Chapman & Hall/CRC Data Mining and Knowledge Discovery Series

RAPIDMINER Data Mining Use Cases and Business Analytics Applications



Edited by Markus Hofmann Ralf Klinkenberg



Contents

1	Introduction to Data Mining and RapidMiner	1
1	What This Book is About and What It is Not Ingo Mierswa	3
	1.1 Introduction 1.2 Coincidence or Not? 1.3 Applications of Data Mining 1.3.1 Financial Services 1.3.2 Retail and Consumer Products 1.3.3 Telecommunications and Media 1.3.4 Manufacturing, Construction, and Electronics 1.4 Fundamental Terms 1.4 Attributes and Target Attributes	$ \begin{array}{r} 3 \\ 4 \\ 7 \\ 7 \\ 8 \\ 9 \\ 10 \\ 11 \\ 11 \\ 11 \end{array} $
	1.4.2Concepts and Examples1.4.3Attribute Roles1.4.4Value Types1.4.5Data and Meta Data1.4.6Modeling	13 14 14 15 16
2	Getting Used to RapidMiner	19
	1 Introduction 2.1 Introduction 2.2 First Start 2.3 Design Perspective 2.4 Building a First Process 2.4.1 Loading Data 2.4.2 Creating a Predictive Model 2.4.3 Executing a Process 2.4.4 Looking at Results	19 19 21 23 24 25 28 29
Π	Basic Classification Use Cases for Credit Approval and in Education	31
3	k-Nearest Neighbor Classification I M. Fareed Akhtar	33
	3.1 Introduction 3.2 Algorithm 3.3 The k-NN Operator in RapidMiner 3.4 Dataset 3.4.1 Teacher Assistant Evaluation Dataset 3.4.2 Basic Information 3.4.3 Examples	33 34 35 35 35 35 35

vii

	3.5 3.6	3.4.4 Operat 3.5.1 3.5.2 3.5.3 3.5.4 3.5.5 3.5.6 3.5.7 3.5.8 Use Ca 3.6.1 3.6.2 3.6.3 3.6.4 3.6.5 2.6.6	Attributes		· · · · · · · · · · · · · · · · · · ·	· · · · · ·	· · · · · · · · · · · · · · · · · · ·	· ·	· · · · · · · · · · · · · · · · · · ·	• • • • • • • • • • • • • •	$35 \\ 36 \\ 36 \\ 36 \\ 37 \\ 37 \\ 37 \\ 38 \\ 38 \\ 38 \\ 39 \\ 39 \\ 40 \\ 40 \\ 41 \\ 41 \\ 41 \\ 41 \\ 41 \\ 41$
	1 17	3.0.0	Model Training, Testing, and Performance Evaluat	LIOI	L	•••	•	• •	•	•	41
4	M = E	arest I areed Al	khtar								45
	4 1	Introd	uction								45
	4.1	Datase	action	• •	•	• •	•	• •	•	•	45
	4.3	Operat	tors Used in This Use Case	• •	•	• •	·	•••	·	·	46
	1.0	431	Read CSV Operator	• •	•	•••	·	• •	·	•	46
		432	Principal Component Analysis Operator	• •	•	• •	•	• •	•	•	47
		433	Split Data Operator	• •	•	• •	•	• •	•	•	48
		434	Performance (Classification) Operator	• •	•	•••	·	•••	·	·	48
	44	Data I	mport	• •	•	•••	·	•••	·	·	48
	4.5	Pre-pr	nessing	• •	•	•••	·	•••	·	·	50
	1.0	4 5 1	Principal Component Analysis	• •	•	•••	·	•••	·	·	50
	46	Model	Training Testing and Performance Evaluation	• •	•	•••	·	•••	·	·	50
	1.0	461	Training the Model	• •	•	•••	·	• •	·	•	51
		462	Testing the Model	• •	•	•••	·	• •	·	•	51
		4.6.3	Performance Evaluation				·				51
-	NT - **	- D									F 9
J	ME	e Daye	b Classification I								JJ
	5 1	Introd	uction								53
	5.2	Datase	action	• •	·	• •	·	• •	·	•	54
	0.2	5.2.1	Credit Approval Dataset	• •	•	• •	·	• •	·	•	54
		5.2.1 5.2.2	Examples	• •	·	• •	·	• •	·	•	54
		5.2.2	Attributes	• •	·	• •	·	• •	·	•	55
	53	0.2.0	tors in This Use Case	• •	•	• •	·	• •	·	·	55
	0.0	531	Rename by Replacing Operator	• •	·	• •	·	• •	·	•	56
		532	Filter Examples Operator	• •	·	• •	·	• •	·	•	56
		5.3.2	Discretize by Binning Operator	•••	•	•••	•	• •	•	•	56
		5.3.3	X-Validation Operator	• •	•	• •	·	• •	•	•	57
		5.0.4 5.3.5	Performance (Binominal Classification) Operator	• •	•	• •	·	• •	·	·	57
	54	Use Co		• •	•	• •	·	• •	·	·	57
	J.T	541	Data Import	•••	•	• •	·	• •	·	·	58
		542	Pre-processing	• •	•	• •	•	• •	•	•	58
		J. 1.4	1 10 Procossing	• •	·	• •	•	• •	•	•	00

5.4.3 Model Training, Testing, and Performance Evaluation 61 M. Farced Akhtar 65 6.1 Dataset 65 6.1.1 Nursery Dataset 65 6.1.2 Basic Information 65 6.1.3 Examples 66 6.4 Attributes 66 6.2 Operators in this Use Case 67 6.2.1 Read Excel Operator 67 6.3.2 Select Attributes Operator 67 6.3.1 Data Import 68 6.3.2 Pre-processing 69 6.3.3 Model Training, Testing, and Performance Evaluation 69 6.3.4 A Deeper Look into the Naïve Bayes Algorithm 71 HI Marketing 77 7.1 Introduction 77 7.2 Business Understanding 78 7.3 Data Understanding 79 7.4 Data Preparation 81 7.4.1 Assembling the Data 82 7.5.2 Class Imbalance 88 7.5.3 Simple Model Evaluation 89 <				ix
6 Naïve Bayes Classificaton II 65 M. Farced Akhtar 65 6.1.1 Nursery Dataset 65 6.1.2 Basic Information 65 6.1.3 Examples 66 6.1.4 Attributes 66 6.2 Operators in this Use Case 67 6.2.1 Read Excel Operator 67 6.2.2 Select Attributes Operator 67 6.3.1 Data Import 68 6.3.2 Pre-processing 69 6.3.3 Model Training, Testing, and Performance Evaluation 69 6.3.4 Deeper Look into the Naïve Bayes Algorithm 71 III Marketing, Cross-Selling, and Recommender System Use Cases Cases 75 7 Who Wants My Product? Affinity-Based Marketing 77 7.2 Business Understanding 79 7.4 Data Understanding 79 7.4 Data Understanding 79 7.4 Data Understanding 81 7.5.1 Continuous Evaluation and Cross Validation 87 7.5.2			5.4.3 Model Training, Testing, and Performance Evaluation	61
M. Fareed Akhtar 6.1 6.1 Dataset 65 6.1.2 Basic Information 65 6.1.3 Examples 66 6.1.4 Attributes 66 6.1.4 Attributes 66 6.2 Operators in this Use Case 67 6.2.1 Read Excel Operator 67 6.3.2 Select Attributes Operator 67 6.3.1 Data Import 68 6.3.2 Pre-processing 69 6.3.3 Model Training, Testing, and Performance Evaluation 69 6.3.4 Deeper Look into the Naïve Bayes Algorithm 71 III Marketing, Cross-Selling, and Recommender System Use Cases Cases 75 Who Wants My Product? Affinity-Based Marketing 77 7.1 Introduction 77 72 7.4 Data Understanding 79 74 7.4 Data Understanding 79 7.4.2 Preparation 81 7.5.1 Continuous Evaluation and Cross Validation 87 7.5.2 Class Imbalance 82	6	Naïv	e Bayes Classificaton II	65
6.1 Dataset 65 6.1.1 Nursery Dataset 65 6.1.2 Basic Information 65 6.1.3 Examples 66 6.1.4 Attributes 66 6.1.4 Attributes 66 6.2 Operators in this Use Case 67 6.2.1 Read Excel Operator 67 6.2.2 Select Attributes Operator 67 6.3.1 Data Import 68 6.3.2 Pre-processing 67 6.3.3 Model Training, Testing, and Performance Evaluation 69 6.3.4 A Deeper Look into the Naïve Bayes Algorithm 71 III Marketing, Cross-Selling, and Recommender System Use Cases 75 7 Who Wants My Product? Affinity-Based Marketing 77 7.1 Introduction 77 7.2 Business Understanding 78 7.3 Data Understanding 78 7.4.1 Assembling the Data 81 7.4.1 Assembling the Data Mining 86 7.5.1 Confidence Values, ROC, and Lift Charts 90 <td></td> <td>M. F</td> <td>areed Akhtar</td> <td></td>		M. F	areed Akhtar	
6.1.1 Nursery Dataset 65 6.1.2 Basic Information 65 6.1.3 Examples 66 6.1.4 Attributes 66 6.1.4 Attributes 66 6.1.4 Attributes 67 6.2.1 Read Excel Operator 67 6.2.2 Select Attributes Operator 67 6.3.1 Data Import 68 6.3.2 Pre-processing 69 6.3.3 Model Training, Testing, and Performance Evaluation 69 6.3.4 Deeper Look into the Naïve Bayes Algorithm 71 III Marketing, Cross-Selling, and Recommender System Use Cases Cases 75 Who Wants My Product? Affinity-Based Marketing 77 7.1 Introduction 77 79 7.4 Data Preparation 78 73 7.5 Data Understanding 79 74.1 Assembling the Data 82 7.4.1 Assembling the Data 82 7.4.2 Preparing for Data Mining 86 7.5.1 Continuous Evaluation and Cross Validation 87 <td></td> <td>6.1</td> <td>Dataset</td> <td>65</td>		6.1	Dataset	65
6.1.2 Basic Information 65 6.1.3 Examples 66 6.1.4 Attributes 66 6.2 Operators in this Use Case 67 6.2.1 Read Excel Operator 67 6.2.2 Select Attributes Operator 67 6.3.1 Data Import 68 6.3.2 Pre-processing 69 6.3.3 Model Training, Testing, and Performance Evaluation 69 6.3.4 A Deeper Look into the Naïve Bayes Algorithm 71 III Marketing, Cross-Selling, and Recommender System Use Cases Cases 75 7 Who Wants My Product? Affinity-Based Marketing 77 7.2 Business Understanding 78 7.3 Data Understanding 78 7.4 Data Preparation 81 7.4.1 Assembling the Data 82 7.5.2 Class Imbalance 88 7.5.3 Simple Model Evaluation 89 7.5.4 Confidence Values, ROC, and Lift Charts 90 7.5.4 Condeling alese Study 97			6.1.1 Nursery Dataset	65
6.1.3 Examples 66 6.1.4 Attributes 66 6.2 Operators in this Use Case 67 6.2.1 Read Excel Operator 67 6.2.2 Select Attributes Operator 67 6.3.1 Data Import 68 6.3.2 Pre-processing 69 6.3.3 Model Training, Testing, and Performance Evaluation 69 6.3.4 Deeper Look into the Naïve Bayes Algorithm 71 III Marketing, Cross-Selling, and Recommender System Use Cases 75 7 Who Wants My Product? Affinity-Based Marketing 77 7.1 Introduction 77 7 7.2 Business Understanding 78 7.3 Data Understanding 78 7.4 Data Preparation 81 7.4.1 Assembling the Data 82 7.4.2 Preparing for Data Mining 86 7.5.1 Continuous Evaluation 87 7.5.2 Class Imbalance 88 7.5.3 Simple Model Evaluation 89 7.5.4 Confidence Values, ROC, and Lif			6.1.2 Basic Information	65
6.1.4 Attributes 66 6.2 Operators in this Use Case 67 6.2.1 Read Excel Operator 67 6.2.2 Select Attributes Operator 67 6.3.1 Data Import 68 6.3.2 Pre-processing 69 6.3.3 Model Training, Testing, and Performance Evaluation 69 6.3.4 A Deeper Look into the Naïve Bayes Algorithm 71 III Marketing, Cross-Selling, and Recommender System Use Cases 75 7 Who Wants My Product? Affinity-Based Marketing 77 7.1 Introduction 77 7.2 Business Understanding 78 7.3 Data Understanding 79 7.4 Data Preparing for Data Mining 86 7.5 Modelling and Evaluation 87 7.5.1 Continuous Evaluation and Cross Validation 87 7.5.2 Class Imbalance 88 7.5.3 Simple Model Evaluation 89 7.5.4 Confidence Values, ROC, and Lift Charts 90 7.5.5 Trying Different Models 92			6.1.3 Examples	66
6.2 Operators in this Use Case 67 6.2.1 Read Excel Operator 67 6.2.2 Select Attributes Operator 67 6.3.1 Data Import 68 6.3.2 Pre-processing 69 6.3.3 Model Training, Testing, and Performance Evaluation 69 6.3.4 A Deeper Look into the Naïve Bayes Algorithm 71 III Marketing, Cross-Selling, and Recommender System Use Cases 75 7 Who Wants My Product? Affinity-Based Marketing 77 7.1 Introduction 77 7.2 Business Understanding 78 7.3 Data Understanding 79 7.4 Data Preparation 81 7.4.1 Assembling the Data 82 7.4.2 Preparing for Data Mining 86 7.5.1 Continuous Evaluation 87 7.5.2 Class Imbalance 88 7.5.3 Simple Model Evaluation 89 7.5.4 Confidence Values, ROC, and Lift Charts 90 7.5.5 Trying Different Models 92 7.6 Depl			6.1.4 Attributes	66
6.2.1 Read Excel Operator 67 6.2.2 Select Attributes Operator 67 6.3.1 Data Import 68 6.3.2 Pre-processing 69 6.3.3 Model Training, Testing, and Performance Evaluation 69 6.3.4 A Deeper Look into the Naïve Bayes Algorithm 71 III Marketing, Cross-Selling, and Recommender System Use Cases 75 7 Who Wants My Product? Affinity-Based Marketing 77 7.1 Introduction 77 7.2 Business Understanding 78 7.3 Data Understanding 78 7.4 Data Preparation 81 7.4.1 Assembling the Data 82 7.4.2 Preparation 87 7.5.1 Continuous Evaluation and Cross Validation 87 7.5.2 Class Imbalance 88 7.5.3 Simple Model Evaluation 89 7.5.4 Confidence Values, ROC, and Lift Charts 90 7.5.5 Trying Different Models 92 7.6 Deployment 93 7.7 Conclusions </td <td></td> <td>6.2</td> <td>Operators in this Use Case</td> <td>67</td>		6.2	Operators in this Use Case	67
6.2.2 Select Attributes Operator 67 6.3 Use Case 67 6.3.1 Data Import 68 6.3.2 Pre-processing 69 6.3.3 Model Training, Testing, and Performance Evaluation 69 6.3.4 A Deeper Look into the Naïve Bayes Algorithm 71 III Marketing, Cross-Selling, and Recommender System Use Cases 75 7 Who Wants My Product? Affinity-Based Marketing 77 <i>Euler Timm</i> 71 71 7.1 Introduction 77 7.2 Business Understanding 78 7.3 Data Understanding 78 7.4 Data Preparation 81 7.4.1 Assembling the Data 82 7.4.2 Preparing for Data Mining 86 7.5.3 Cast Imbalance 88 7.5.4 Continuous Evaluation 87 7.5 Class Imbalance 90 7.5.4 Confidence Values, ROC, and Lift Charts 90 7.5.5 Trying Different Models 92 7.6 Deployment 93			6.2.1 Read Excel Operator	67
6.3 Use Case 67 6.3.1 Data Import 68 6.3.2 Pre-processing 69 6.3.3 Model Training, Testing, and Performance Evaluation 69 6.3.4 A Deeper Look into the Naïve Bayes Algorithm 71 III Marketing, Cross-Selling, and Recommender System Use Cases 75 7 Who Wants My Product? Affinity-Based Marketing 77 7.1 Introduction 77 7.2 Business Understanding 78 7.3 Data Understanding 79 7.4 Data Preparation 81 7.4.1 Assembling the Data 82 7.4.2 Preparing for Data Mining 86 7.5 Modelling and Evaluation 87 7.5.1 Continuous Evaluation and Cross Validation 87 7.5.2 Class Imbalance 88 7.5.3 Simple Model Evaluation 89 7.5.4 Confidence Values, ROC, and Lift Charts 90 7.5.5 Trying Different Models 92 7.6 Deployment 93 7.7 Conclusions			6.2.2 Select Attributes Operator	67
6.3.1 Data Import 68 6.3.2 Pre-processing 69 6.3.3 Model Training, Testing, and Performance Evaluation 69 6.3.4 A Deeper Look into the Naïve Bayes Algorithm 71 III Marketing, Cross-Selling, and Recommender System Use Cases 75 7 Who Wants My Product? Affinity-Based Marketing 77 <i>Euler Timm</i> 71 71 7.1 Introduction 77 7.2 Business Understanding 78 7.3 Data Understanding 79 7.4 Data Preparation 81 7.4.1 Assembling the Data 82 7.4.2 Preparing for Data Mining 86 7.5.4 Confidence Valuation and Cross Validation 87 7.5.2 Class Imbalance 88 7.5.3 Simple Model Evaluation 89 7.5.4 Confidence Values, ROC, and Lift Charts 90 7.5.5 Trying Different Models 92 7.6 Deployment 93 7.7 Conclusions 94 8 Basic Association Rule Mining		6.3	Use Case	67
6.3.2 Pre-processing			6.3.1 Data Import	68
6.3.3 Model Training, Testing, and Performance Evaluation 69 6.3.4 A Deeper Look into the Naïve Bayes Algorithm 71 III Marketing, Cross-Selling, and Recommender System Use Cases 75 7 Who Wants My Product? Affinity-Based Marketing 77 Euler Timm 71 7.1 Introduction 77 7.2 Business Understanding 78 7.3 Data Understanding 79 7.4 Data Preparation 81 7.4.1 Assembling the Data 82 7.4.2 Preparing for Data Mining 86 7.5 Model Evaluation 87 7.5.1 Continuous Evaluation and Cross Validation 87 7.5.2 Class Imbalance 88 7.5.3 Simple Model Evaluation 89 7.5.4 Confidence Values, ROC, and Lift Charts 90 7.5.5 Trying Different Models 92 7.6 Deployment 93 7.7 Conclusions 94 8 Basic Association Rule Mining in RapidMiner 97 9.4 Basic Association			6.3.2 Pre-processing	69
6.3.4 A Deeper Look into the Naïve Bayes Algorithm 71 III Marketing, Cross-Selling, and Recommender System Use Cases 75 7 Who Wants My Product? Affinity-Based Marketing 77 7.1 Introduction 77 7.2 Business Understanding 78 7.3 Data Understanding 79 7.4 Data Preparation 81 7.4.1 Assembling the Data 82 7.4.2 Preparing for Data Mining 86 7.5 Modelling and Evaluation 87 7.5.2 Class Imbalance 88 7.5.3 Simple Model Evaluation 89 7.5.4 Confidence Values, ROC, and Lift Charts 90 7.5.5 Trying Different Models 92 7.6 Deployment 93 7.7 Conclusions 94 8 Basic Association Rule Mining in RapidMiner 97 91 Introduction 119 Matej Mihelčić, Matko Bošnjak, Nino Antulov-Fantulin, and Tomislav Šmuc 91 9.1 Introduction 120 9.2 The Recommender Extension 121 9.2.1 Recommendation Operators 121 9.2.2 Data Format 122 9.3 The VideoLectures.net Dataset 126			6.3.3 Model Training, Testing, and Performance Evaluation	69
III Marketing, Cross-Selling, and Recommender System Use Cases 75 7 Who Wants My Product? Affinity-Based Marketing 77 <i>Euler Timm</i> 77 7.1 Introduction 77 7.2 Business Understanding 78 7.3 Data Understanding 79 7.4 Data Preparation 81 7.4.1 Assembling the Data 82 7.4.2 Preparing for Data Mining 86 7.5 Modelling and Evaluation 87 7.5.2 Class Imbalance 89 7.5.3 Simple Model Evaluation 89 7.5.4 Confidence Values, ROC, and Lift Charts 90 7.5.5 Trying Different Models 92 7.6 Deployment 93 7.7 Conclusions 94 8 Basic Association Rule Mining in RapidMiner 97 94 Basic Association Rule Mining in RapidMiner 119 Matej Mihelčić, Matko Bošnjak, Nino Antulov-Fantulin, and Tomislav Šmuc 91 9.1 Introduction 120 9.2 The Recommender Extension 12			6.3.4 A Deeper Look into the Naïve Bayes Algorithm	71
IIIWith Recting, Oross-Schning, and Tecconfinented System OstCases75 7 Who Wants My Product? Affinity-Based Marketing77Euler Timm777.1Introduction777.2Business Understanding797.4Data Understanding797.4Data Understanding797.4Data Preparation817.4.1Assembling the Data827.4.2Preparing for Data Mining867.5.4Continuous Evaluation877.5.2Class Imbalance887.5.3Simple Model Evaluation897.5.4Confidence Values, ROC, and Lift Charts907.5.5Trying Different Models927.6Deployment937.7Conclusions948Basic Association Rule Mining in RapidMiner979Constructing Recommender Systems in RapidMiner119Mattej Mihelčić, Matko Bošnjak, Nino Antulov-Fantulin, and Tomislav Šmuc919.1Introduction1219.2.1Recommender Extension1219.2.2Data Format1229.3Preformance Measures1249.3The VideoLectures net Dataset1269.4Collaborative-based Systems127	тт	τī	Jarketing Cross-Selling and Recommender System Use	
7Who Wants My Product? Affinity-Based Marketing77Euler Timm71Introduction772Business Understanding787.1Introduction777.2Business Understanding797.4Data Preparation817.4.1Assembling the Data827.4.2Preparing for Data Mining867.5Modelling and Evaluation877.5.1Continuous Evaluation and Cross Validation877.5.2Class Imbalance887.5.3Simple Model Evaluation897.5.4Confidence Values, ROC, and Lift Charts907.5.5Trying Different Models927.6Deployment937.7Conclusions948Basic Association Rule Mining in RapidMiner979Constructing Recommender Systems in RapidMiner19Matej Mihelčić, Matko Bošnjak, Nino Antulov-Fantulin, and Tomislav Šmuc9.19.1Introduction1209.2Data Format1229.2.3Performance Measures1249.3The VideoLectures.net Dataset1269.4Collaborative-based Systems1269.4Collaborative-based Systems127	11	1 N (Cases	75
7 Who Wants My Product? Affinity-Based Marketing 77 Euler Timm 71 7.1 Introduction 77 7.2 Business Understanding 78 7.3 Data Understanding 79 7.4 Data Preparation 81 7.4.1 Assembling the Data 82 7.4.2 Preparing for Data Mining 86 7.5 Modelling and Evaluation 87 7.5.1 Continuous Evaluation and Cross Validation 87 7.5.2 Class Imbalance 88 7.5.3 Simple Model Evaluation 89 7.5.4 Confidence Values, ROC, and Lift Charts 90 7.5.5 Trying Different Models 92 7.6 Deployment 93 7.7 Conclusions 94 8 Basic Association Rule Mining in RapidMiner 97 9 Constructing Recommender Systems in RapidMiner 119 Matej Mihelčić, Matko Bošnjak, Nino Antulov-Fantulin, and Tomislav Šmuc 91 9.1 Introduction 120 9.2 The Recommender Extension 121 <th></th> <th></th> <th></th> <th></th>				
7.1Introduction777.2Business Understanding787.3Data Understanding797.4Data Preparation817.4.1Assembling the Data827.4.2Preparing for Data Mining867.5Modelling and Evaluation877.5.1Continuous Evaluation and Cross Validation877.5.2Class Imbalance887.5.3Simple Model Evaluation897.5.4Confidence Values, ROC, and Lift Charts907.5.5Trying Different Models927.6Deployment937.7Conclusions948Basic Association Rule Mining in RapidMiner979Constructing Recommender Systems in RapidMiner119Matej Mihelčić, Matko Bošnjak, Nino Antulov-Fantulin, and Tomislav Šmuc1209.2The Recommender Extension1219.2.1Recommender Extension1219.2.2Data Format1229.2.3Performance Measures1249.3The VideoLectures.net Dataset1269.4Collaborative-based Systems127	7	Who Euler	• Wants My Product? Affinity-Based Marketing • Timm	77
7.2 Business Understanding 78 7.3 Data Understanding 79 7.4 Data Preparation 81 7.4.1 Assembling the Data 82 7.4.2 Preparing for Data Mining 86 7.5 Modelling and Evaluation 87 7.5.1 Continuous Evaluation and Cross Validation 87 7.5.2 Class Imbalance 88 7.5.3 Simple Model Evaluation 89 7.5.4 Confidence Values, ROC, and Lift Charts 90 7.5.5 Trying Different Models 92 7.6 Deployment 93 7.7 Conclusions 94 8 Basic Association Rule Mining in RapidMiner 97 Matthew A. North 8.1 Data Mining Case Study 97 9 Constructing Recommender Systems in RapidMiner 119 Matej Mihelčić, Matko Bošnjak, Nino Antulov-Fantulin, and Tomislav Šmuc 9.1 9.1 Introduction 120 9.2 The Recommender Extension 121 9.2.1 Recommender Extension 121 9.2.2 <		7.1	Introduction	77
7.3Data Understanding7.97.4Data Preparation817.4.1Assembling the Data827.4.2Preparing for Data Mining867.5Modelling and Evaluation877.5.1Continuous Evaluation and Cross Validation877.5.2Class Imbalance887.5.3Simple Model Evaluation897.5.4Confidence Values, ROC, and Lift Charts907.5.5Trying Different Models927.6Deployment937.7Conclusions948Basic Association Rule Mining in RapidMiner97Matthew A. North8.1Data Mining Case Study979Constructing Recommender Systems in RapidMiner119Matej Mihelčić, Matko Bošnjak, Nino Antulov-Fantulin, and Tomislav Šmuc9.19.1Introduction1209.2The Recommender Extension1219.2.1Recommender Extension1219.2.2Data Format1229.3The VideoLectures.net Dataset1249.3The VideoLectures.de Systems127		7.2	Business Understanding	78
7.4Data Preparation817.4.1Assembling the Data827.4.2Preparing for Data Mining867.5.1Continuous Evaluation877.5.2Class Imbalance887.5.3Simple Model Evaluation897.5.4Confidence Values, ROC, and Lift Charts907.5.5Trying Different Models927.6Deployment937.7Conclusions948Basic Association Rule Mining in RapidMiner97Matthew A. North8.1Data Mining Case Study979Constructing Recommender Systems in RapidMiner119Matej Mihelčić, Matko Bošnjak, Nino Antulov-Fantulin, and Tomislav Šmuc9.19.1Introduction1209.2The Recommender Extension1219.2.1Recommender Extension1219.2.2Data Format1229.2.3Performance Measures1249.3The VideoLectures.net Dataset1269.4Collaborative-based Systems127		7.3	Data Understanding	79
7.4.1 Assembling the Data 82 7.4.2 Preparing for Data Mining 86 7.5 Modelling and Evaluation 87 7.5.1 Continuous Evaluation and Cross Validation 87 7.5.2 Class Imbalance 88 7.5.3 Simple Model Evaluation 89 7.5.4 Confidence Values, ROC, and Lift Charts 90 7.5.5 Trying Different Models 92 7.6 Deployment 93 7.7 Conclusions 94 8 Basic Association Rule Mining in RapidMiner 97 Matthew A. North 97 9.1 Introduction 119 Matej Mihelčić, Matko Bošnjak, Nino Antulov-Fantulin, and Tomislav Šmuc 91 9.1 Introduction 120 9.2 The Recommender Extension 121 9.2.1 Recommendation Operators 121 9.2.2 Data Format 122 9.2.3 Performance Measures 124 9.3 The VideoLectures.net Dataset 126 9.4 Collaborative-based Systems 127 </td <td></td> <td>7.4</td> <td>Data Preparation</td> <td>81</td>		7.4	Data Preparation	81
7.4.2 Preparing for Data Mining 86 7.5 Modelling and Evaluation 87 7.5.1 Continuous Evaluation and Cross Validation 87 7.5.2 Class Imbalance 88 7.5.3 Simple Model Evaluation 89 7.5.4 Confidence Values, ROC, and Lift Charts 90 7.5.5 Trying Different Models 92 7.6 Deployment 93 7.7 Conclusions 94 8 Basic Association Rule Mining in RapidMiner 97 Matthew A. North 97 9 Constructing Recommender Systems in RapidMiner 119 Matej Mihelčić, Matko Bošnjak, Nino Antulov-Fantulin, and Tomislav Šmuc 9.1 9.1 Introduction 120 9.2 The Recommender Extension 121 9.2.1 Recommender Extension 121 9.2.2 Data Format 122 9.2.3 Performance Measures 124 9.3 The VideoLectures.net Dataset 126 9.4 Collaborative-based Systems 127			7.4.1 Assembling the Data	82
7.5 Modelling and Evaluation 87 7.5.1 Continuous Evaluation and Cross Validation 87 7.5.2 Class Imbalance 88 7.5.3 Simple Model Evaluation 89 7.5.4 Confidence Values, ROC, and Lift Charts 90 7.5.5 Trying Different Models 92 7.6 Deployment 93 7.7 Conclusions 93 7.7 Conclusions 94 8 Basic Association Rule Mining in RapidMiner 97 Matthew A. North 97 9 Constructing Recommender Systems in RapidMiner 119 Matej Mihelčić, Matko Bošnjak, Nino Antulov-Fantulin, and Tomislav Šmuc 120 9.1 Introduction 120 9.2 The Recommender Extension 121 9.2.1 Recommender Extension 122 9.2.3 Performance Measures 124 9.3 The VideoLectures.net Dataset 126 9.4 Collaborative-based Systems 127			7.4.2 Preparing for Data Mining	86
7.5.1 Continuous Evaluation and Cross Validation 87 7.5.2 Class Imbalance 88 7.5.3 Simple Model Evaluation 89 7.5.4 Confidence Values, ROC, and Lift Charts 90 7.5.5 Trying Different Models 92 7.6 Deployment 93 7.7 Conclusions 94 8 Basic Association Rule Mining in RapidMiner 97 Matthew A. North 91 97 9 Constructing Recommender Systems in RapidMiner 97 9 Constructing Recommender Systems in RapidMiner 119 Matej Mihelčić, Matko Bošnjak, Nino Antulov-Fantulin, and Tomislav Šmuc 9.1 9.1 Introduction 120 9.2 The Recommender Extension 121 9.2.1 Recommender Extension 121 9.2.2 Data Format 122 9.2.3 Performance Measures 124 9.3 The VideoLectures.net Dataset 126 9.4 Collaborative-based Systems 127		7.5	Modelling and Evaluation	87
7.5.2 Class Imbalance 88 7.5.3 Simple Model Evaluation 89 7.5.4 Confidence Values, ROC, and Lift Charts 90 7.5.5 Trying Different Models 92 7.6 Deployment 93 7.7 Conclusions 94 8 Basic Association Rule Mining in RapidMiner 97 Matthew A. North 97 9 Constructing Recommender Systems in RapidMiner 97 9 Constructing Recommender Systems in RapidMiner 119 Matej Mihelčić, Matko Bošnjak, Nino Antulov-Fantulin, and Tomislav Šmuc 9.1 9.1 Introduction 120 9.2 The Recommender Extension 121 9.2.1 Recommender Extension 121 9.2.2 Data Format 122 9.2.3 Performance Measures 124 9.3 The VideoLectures.net Dataset 126 9.4 Collaborative-based Systems 127			7.5.1 Continuous Evaluation and Cross Validation	87
7.5.3 Simple Model Evaluation 89 7.5.4 Confidence Values, ROC, and Lift Charts 90 7.5.5 Trying Different Models 92 7.6 Deployment 93 7.7 Conclusions 94 8 Basic Association Rule Mining in RapidMiner 97 Matthew A. North 97 9 Constructing Recommender Systems in RapidMiner 97 9 Constructing Recommender Systems in RapidMiner 119 Matej Mihelčić, Matko Bošnjak, Nino Antulov-Fantulin, and Tomislav Šmuc 91 9.1 Introduction 120 9.2 The Recommender Extension 121 9.2.1 Recommender Extension 121 9.2.2 Data Format 122 9.3 The VideoLectures.net Dataset 124 9.3 The VideoLectures.net Dataset 126 9.4 Collaborative-based Systems 127			7.5.2 Class Imbalance	88
7.5.4 Confidence Values, ROC, and Lift Charts 90 7.5.5 Trying Different Models 92 7.6 Deployment 93 7.7 Conclusions 94 8 Basic Association Rule Mining in RapidMiner 97 Matthew A. North 97 9 Constructing Recommender Systems in RapidMiner 97 9 Constructing Recommender Systems in RapidMiner 119 Matej Mihelčić, Matko Bošnjak, Nino Antulov-Fantulin, and Tomislav Šmuc 91 9.1 Introduction 120 9.2 The Recommender Extension 121 9.2.1 Recommendation Operators 121 9.2.2 Data Format 122 9.3.3 Performance Measures 124 9.3 The VideoLectures.net Dataset 126 9.4 Collaborative-based Systems 127			7.5.3 Simple Model Evaluation	89
7.5.5 Trying Different Models 92 7.6 Deployment 93 7.7 Conclusions 94 8 Basic Association Rule Mining in RapidMiner 97 Matthew A. North 97 9 Constructing Recommender Systems in RapidMiner 119 Matej Mihelčić, Matko Bošnjak, Nino Antulov-Fantulin, and Tomislav Šmuc 120 9.1 Introduction 121 9.2.1 Recommender Extension 121 9.2.2 Data Format 122 9.2.3 Performance Measures 124 9.3 The VideoLectures.net Dataset 126 9.4 Collaborative-based Systems 127			7.5.4 Confidence Values, ROC, and Lift Charts	90
7.6 Deployment 93 7.7 Conclusions 94 8 Basic Association Rule Mining in RapidMiner 97 Matthew A. North 97 9 Constructing Recommender Systems in RapidMiner 97 9 Constructing Recommender Systems in RapidMiner 119 Matej Mihelčić, Matko Bošnjak, Nino Antulov-Fantulin, and Tomislav Šmuc 120 9.1 Introduction 120 9.2 The Recommender Extension 121 9.2.1 Recommendation Operators 121 9.2.2 Data Format 122 9.3.3 Performance Measures 124 9.3 The VideoLectures.net Dataset 126 9.4 Collaborative-based Systems 127			7.5.5 Trying Different Models	92
7.7 Conclusions 94 8 Basic Association Rule Mining in RapidMiner 97 Matthew A. North 97 8.1 Data Mining Case Study 97 9 Constructing Recommender Systems in RapidMiner 97 9 Constructing Recommender Systems in RapidMiner 119 Matej Mihelčić, Matko Bošnjak, Nino Antulov-Fantulin, and Tomislav Šmuc 9.1 9.1 Introduction 120 9.2 The Recommender Extension 121 9.2.1 Recommendation Operators 121 9.2.2 Data Format 122 9.2.3 Performance Measures 124 9.3 The VideoLectures.net Dataset 126 9.4 Collaborative-based Systems 127		7.6	Deployment	93
8 Basic Association Rule Mining in RapidMiner 97 Matthew A. North 97 8.1 Data Mining Case Study 97 9 Constructing Recommender Systems in RapidMiner 119 Matej Mihelčić, Matko Bošnjak, Nino Antulov-Fantulin, and Tomislav Šmuc 119 9.1 Introduction 120 9.2 The Recommender Extension 121 9.2.1 Recommendation Operators 121 9.2.2 Data Format 122 9.2.3 Performance Measures 124 9.3 The VideoLectures.net Dataset 126 9.4 Collaborative-based Systems 127		7.7	Conclusions	94
Matthew A. North 97 8.1 Data Mining Case Study 97 9 Constructing Recommender Systems in RapidMiner 119 Matej Mihelčić, Matko Bošnjak, Nino Antulov-Fantulin, and Tomislav Šmuc 119 9.1 Introduction 120 9.2 The Recommender Extension 121 9.2.1 Recommendation Operators 121 9.2.2 Data Format 122 9.2.3 Performance Measures 124 9.3 The VideoLectures.net Dataset 126 9.4 Collaborative-based Systems 127	8	Basi	c Association Rule Mining in RapidMiner	97
8.1 Data Mining Case Study 97 9 Constructing Recommender Systems in RapidMiner 119 Matej Mihelčić, Matko Bošnjak, Nino Antulov-Fantulin, and Tomislav Šmuc 119 9.1 Introduction 120 9.2 The Recommender Extension 121 9.2.1 Recommendation Operators 121 9.2.2 Data Format 122 9.2.3 Performance Measures 124 9.3 The VideoLectures.net Dataset 126 9.4 Collaborative-based Systems 127	0	Mattl	hew A. North	01
9 Constructing Recommender Systems in RapidMiner 119 Matej Mihelčić, Matko Bošnjak, Nino Antulov-Fantulin, and Tomislav Šmuc 120 9.1 Introduction 120 9.2 The Recommender Extension 121 9.2.1 Recommendation Operators 121 9.2.2 Data Format 122 9.3 The VideoLectures.net Dataset 124 9.3 Collaborative-based Systems 126		8.1	Data Mining Case Study	97
Matej Mihelčić, Matko Bošnjak, Nino Antulov-Fantulin, and Tomislav Šmuc 9.1 Introduction 120 9.2 The Recommender Extension 121 9.2.1 Recommendation Operators 121 9.2.2 Data Format 122 9.2.3 Performance Measures 124 9.3 The VideoLectures.net Dataset 126 9.4 Collaborative-based Systems 127	9	Cons	structing Recommender Systems in RapidMiner	19
9.1 Introduction 120 9.2 The Recommender Extension 121 9.2.1 Recommendation Operators 121 9.2.2 Data Format 122 9.2.3 Performance Measures 124 9.3 The VideoLectures.net Dataset 126 9.4 Collaborative-based Systems 127		Mate	i Mihelčić. Matko Bošniak. Nino Antulov-Fantulin. and Tomislav Šmuc	
9.2 The Recommender Extension 121 9.2.1 Recommendation Operators 121 9.2.2 Data Format 121 9.2.3 Performance Measures 122 9.3 The VideoLectures.net Dataset 126 9.4 Collaborative-based Systems 127		9.1	Introduction	120
9.2.1 Recommendation Operators 121 9.2.2 Data Format 122 9.2.3 Performance Measures 124 9.3 The VideoLectures.net Dataset 126 9.4 Collaborative-based Systems 127		9.2	The Recommender Extension	121
9.2.2 Data Format 122 9.2.3 Performance Measures 124 9.3 The VideoLectures.net Dataset 126 9.4 Collaborative-based Systems 127			9.2.1 Recommendation Operators	121
9.2.3 Performance Measures 124 9.3 The VideoLectures.net Dataset 126 9.4 Collaborative-based Systems 127			9.2.2 Data Format	122
9.3 The VideoLectures.net Dataset 126 9.4 Collaborative-based Systems 127			9.2.3 Performance Measures	124
9.4 Collaborative-based Systems		9.3	The VideoLectures.net Dataset	126
		9.4	Collaborative-based Systems	127

	9.4.1 Neighbourhood-based Recommender Systems 9.4.2 Factorization-based Recommender Systems 9.4.3 Collaborative Recommender Workflows	$127 \\ 128 \\ 130$
	9.4.4 Iterative Online Updates	131
9.5	Content-based Recommendation	132
0.0	9.5.1 Attribute-based Content Recommendation	133
	9.5.2 Similarity-based Content Recommendation	134
9.6	Hybrid Recommender Systems	135
9.7	Providing RapidMiner Recommender System Workflows as Web Services	100
0.1	Using Rapid Analytics	138
	9.7.1 Simple Recommender System Web Service	138
	9.7.2 Guidelines for Optimizing Workflows for Service Usage	139
9.8	Summary	141
10 Reco	ommender System for Selection of the Right Study Program for	•
High	ner Education Students	145
Mila	n Vukićević, Miloš Jovanović, Boris Delibašić, and Milija Suknović	
10.1	Introduction	146
10.2	Literature Review	146
10.3	Automatic Classification of Students using RapidMiner	147
	10.3.1 Data	147
	10.3.2 Processes	147
	10.3.2.1 Simple Evaluation Process	150
	10.3.2.2 Complex Process (with Feature Selection)	152
10.4	Results	154
10.5	Conclusion	155
IV C	Clustering in Medical and Educational Domains	157
		1 50
	alising Clustering Validity Measures	159
Andr	ew Chisholm	1.00
11.1	Overview	160
11.2		160
11.0	11.2.1 A Brief Explanation of k-Means.	101
11.3	Cluster Validity Measures	101
	11.3.1 Internal validity Measures $\dots \dots \dots$	101
	11.3.2 External Validity Measures	162
11.4	11.3.3 Relative Validity Measures	163
11.4		163
	11.4.1 Artificial Data	164
11 5	11.4.2 E-coli Data	164
11.5	Setup	165
	11.5.1 Download and Install R Extension	166
11.0	11.5.2 Processes and Data	166
11.6	The Process in Detail	167
	11.6.1 Import Data (A) \ldots (D)	168
	11.6.2 Generate Clusters (B) $\ldots \ldots \ldots$	169
	11.6.3 Generate Ground Truth Validity Measures (C)	170
	11.6.4 Generate External Validity Measures (D)	172
	11.6.5 Generate Internal Validity Measures (E)	173
	11.6.6 Output Results (F)	174

х

11.7 11.8	Running the Process and Displaying Results	175 176 176 176 178
11.9	Conclusion	. 181
12 Grou Milan	uping Higher Education Students with RapidMiner n Vukićević, Miloš Jovanović, Boris Delibašić, and Milija Suknović	185
12.1	Introduction	. 185
12.2	Related Work	. 186
12.3	Using RapidMiner for Clustering Higher Education Students	. 186
	12.3.1 Data	. 187
	12.3.2 Process for Automatic Evaluation of Clustering Algorithms	. 187
	12.3.3 Results and Discussion	. 191
12.4	Conclusion	. 193
V T C	fext Mining: Spam Detection, Language Detection, a Customer Feedback Analysis	nd 197
13 Dete	ecting Text Message Spam	199
Neil	McGuigan	200
13.1		. 200
10.2 12.2	Apprying This Technique in Other Domains	200
10.0	Catting the Data	200
15.4 12 5	Getting the Data	201
15.0	Loading the Text	201
	13.5.1 Data Import Wizard Step 1	201
	13.5.2 Data Import Wizard Step 2	202
	13.5.5 Data Import Wizard Step 5	202
	13.5.4 Data import wizard Step 4	202
126	Framining the Text	202
13.0	12.6.1 Tokonizing the Decument	· · 205
	13.6.2 Creating the Word List and Word Voctor	203
	13.6.2 Examining the Word Vector	204
137	Processing the Text for Classification	204
10.1	13.7.1 Text Processing Concents	206
13.8	The Naïve Bayes Algorithm	200
10.0	13.8.1 How It Works	207
13.9	Classifying the Data as Spam or Ham	208
13.10	Validating the Model	208
13.10	Applying the Model to New Data	209
10.11	13 11 1 Bunning the Model on New Data	210
13 12	2 Improvements	210
13 13	Summary	211
10.10	·	• • • • • •
14 Rob	ust Language Identification with RapidMiner: A Text Mining V	Use
Case		213
Matk	o Bosnjak, Eduarda Mendes Rodrigues, and Luis Sarmento	01.4
14.1	Introduction	. 214
14.2	The Froblem of Language Identification	. 215

xi

	14.3	Text Representation	17
		14.3.1 Encoding	17
		14.3.2 Token-based Representation	18
		14.3.3 Character-Based Representation	19
		14.3.4 Bag-of-Words Representation	19
	14.4	Classification Models	20
	14.5	Implementation in RapidMiner	21
		14.5.1 Datasets	21
		14.5.2 Importing Data	23
		14.5.3 Frequent Words Model	25
		14.5.4 Character n-Grams Model	29
		14.5.5 Similarity-based Approach	32
	14.6	Application	34
		14.6.1 RapidAnalytics	34
		14.6.2 Web Page Language Identification	34
	14.7	Summary 25	37
15	\mathbf{Text}	Mining with RapidMiner 24	11
	Gurda	al Ertek, Dilek Tapucu, and Inanc Arin	
	15.1	Introduction	42
		15.1.1 Text Mining	42
		15.1.2 Data Description	42
		15.1.3 Running RapidMiner	42
		15.1.4 RapidMiner Text Processing Extension Package	43
		15.1.5 Installing Text Mining Extensions	43
	15.2	Association Mining of Text Document Collection (Process01) 24	43
		15.2.1 Importing Process 01	43
		15.2.2 Operators in Process01	43
		15.2.3 Saving Process01	17
	15.3	Clustering Text Documents (Process02) 24	48
		15.3.1 Importing Process02	48
		15.3.2 Operators in Process02	48
		15.3.3 Saving Process02	50
	15.4	Running Process01 and Analyzing the Results	50
		15.4.1 Running Process01	50
		15.4.2 Empty Results for Process01	52
		15.4.3 Specifying the Source Data for Process01	52
		15.4.4 Re-Running Process01	53
		15.4.5 Process01 Results	53
		15.4.6 Saving Process01 Results	57
	15.5	Running Process02 and Analyzing the Results	57
		15.5.1 Running Process02	57
		15.5.2 Specifying the Source Data for Process02	57
		15.5.3 Process02 Results	57

xii

15.6	Conclusions	261
	Feature Selection and Classification in Astroparticle Physics and in Medical Domains	e 263
16 App Tim	Dication of RapidMiner in Neutrino Astronomy	265
16.1	Protons Photons and Neutrinos	265
16.2	Neutrino Astronomy	$200 \\ 267$
16.2	Feature Selection	260
10.0	16.3.1 Installation of the Feature Selection Extension for BanidMiner	269
	16.3.2 Feature Selection Setup	$\frac{203}{270}$
	16.3.3 Inner Process of the LOOP PARAMETERS Operator	$\frac{210}{271}$
	16.3.4 Inner Operators of the WRAPPER X-VALIDATION	271
	16.3.5 Settings of the LOOP PARAMETERS Operator	274
	16.3.6 Feature Selection Stability	275
16.4	Event Selection Using a Bandom Forest	$210 \\ 277$
10.1	16.4.1 The Training Setup	278
	16.4.2 The Bandom Forest in Greater Detail	280
	16.4.3 The Bandom Forest Settings	281
	16.4.4 The Testing Setup	282
16.5	Summary and Outlook	$\frac{282}{285}$
17 Mec	dical Data Mining	289
Mer	tik Matei and Palfu Miroslav	200
17 1	Background	290
17.1 17.2	Description of Problem Domain: Two Medical Examples	291
± =	17.2.1 Carpal Tunnel Syndrome	291
	17.2.2 Diabetes	292
17.3	Data Mining Algorithms in Medicine	292
1110	17.3.1 Predictive Data Mining	292
	17.3.2 Descriptive Data Mining	293
	17.3.3 Data Mining and Statistics: Hypothesis Testing	294
17.4	Knowledge Discovery Process in BapidMiner: Carpal Tunnel Syndrome	295
	17.4.1 Defining the Problem. Setting the Goals	295
	17.4.2 Dataset Representation	295
	17.4.3 Data Preparation	296
	17.4.4 Modeling	298
	17.4.5 Selecting Appropriate Methods for Classification	298
	17.4.6 Results and Data Visualisation	303
	17.4.7 Interpretation of the Results	303
	17.4.8 Hypothesis Testing and Statistical Analysis	304
	17.4.9 Results and Visualisation	308
17.5	Knowledge Discovery Process in RapidMiner: Diabetes	308
	17.5.1 Problem Definition, Setting the Goals	309
	17.5.2 Data Preparation	309
	17.5.3 Modeling	310
	17.5.4 Results and Data Visualization	312
	17.5.5 Hypothesis Testing	313
17.6	Specifics in Medical Data Mining	314
17.7	Summary	314

xiii

VII		Iolecular Structure- and Property-Activity Relationship Iodeling in Biochemistry and Medicine) 319
18 U N	Using Mode	g PaDEL to Calculate Molecular Properties and Chemoinformatic	c 321
Λ	Mark	us Muehlbacher and Johannes Kornhuber	
1	8.1	Introduction	321
1	8.2	Molecular Structure Formats for Chemoinformatics	321
1	8.3	Installation of the PaDEL Extension for RapidMiner	322
1	8.4	Applications and Capabilities of the PaDEL Extension	323
1	8.5	Examples of Computer-aided Predictions	324
1	8.6	Calculation of Molecular Properties	325
1	8.7	Generation of a Linear Regression Model	325
1	8.8	Example Workflow	326
1	8.9	Summary	328
19 C	Chen	noinformatics: Structure- and Property-activity Relationship Devel	-
0	pme	ent	331
Λ	Mark	us Muehlbacher and Johannes Kornhuber	
1	9.1	Introduction	331
1	9.2	Example Workflow	332
1	9.3	Importing the Example Set	332
1	9.4	Preprocessing of the Data	333
1	9.5	Feature Selection	334
1	9.6	Model Generation	335
1	9.7	Validation	337
1	9.8	Y-Randomization	338
1	9.9	Results	339
1	9.10	Conclusion/Summary	340
VII	I I	mage Mining: Feature Extraction, Segmentation, and	1
	C	Classification	345
20 I	mag Radin	e Mining Extension for RapidMiner (Introductory) n Burget Václay Uber, and Jan Masek	347
2	10000	Introduction	348
2	20.1	Image Reading/Writing	349
2	0.2	Conversion between Colour and Gravscale Images	352
2	20.0	Feature Extraction	353
2	10.1	20.4.1 Local Level Feature Extraction	354
		20.4.2 Segment-Level Feature Extraction	356
		20.4.3 Global-Level Feature Extraction	358
2	20.5	Summary	359
01 T			000
41 I T	mag	e winning Extension for Rapidwiner (Advanced)	303
1	v acia	v Uner and Kaaim Burget	969
2	1.1 1 0		303 964
2	1.2	111 age Classification	304 264
		21.2.1 LOAU IIIIages alle Assigli Labels	304 265
0	01-9	21.2.2 GIODAI FEATURE EXTRACTION	000 960
2	6.1.		308

 xiv

	$21.4 \\ 21.5$	21.3.1 Process Creation Image Segmentation and Feature Extraction Summary	370 372 373
D	K A C	Anomaly Detection, Instance Selection, and Prototype Construction	375
22	Insta Marc	ance Selection in RapidMiner in Blachnik and Miroslaw Kordos	377
	22.1	Introduction	377
	22.2	Instance Selection and Prototype-Based Rule Extension	378
	22.3	Instance Selection	379
		22.3.1 Description of the Implemented Algorithms	381
		22.3.2 Accelerating I-NN Classification	384
		22.3.3 Outlier Elimination and Noise Reduction	389
	00.4	22.3.4 Advances in Instance Selection	392
	22.4 22.5	Mining Lange Detector	395
	22.0 22.6	Summary	401
	22.0		400
23	Anoi Mark	maly Detection	409
	23 1	Introduction	410
	20.1 23.2	Categorizing an Anomaly Detection Problem	412
	20.2	23.2.1 Type of Anomaly Detection Problem (Pre-processing)	412
		23.2.2 Local versus Global Problems	416
		23.2.3 Availability of Labels	416
	23.3	A Simple Artificial Unsupervised Anomaly Detection Example	417
	23.4	Unsupervised Anomaly Detection Algorithms	419
		23.4.1 k-NN Global Anomaly Score	419
		23.4.2 Local Outlier Factor (LOF)	420
		23.4.3 Connectivity-Based Outlier Factor (COF)	421
		23.4.4 Influenced Outlierness (INFLO)	422
		23.4.5 Local Outlier Probability (LoOP)	422
		23.4.6 Local Correlation Integral (LOCI) and aLOCI	422
		23.4.7 Cluster-Based Local Outlier Factor (CBLOF)	423
		23.4.8 Local Density Cluster-Based Outlier Factor (LDCOF)	424
	23.5	An Advanced Unsupervised Anomaly Detection Example	425
	23.6	Semi-supervised anomaly detection	428
		23.6.1 Using a One-Class Support Vector Machine (SVM)	428
		23.6.2 Clustering and distance computations for detecting anomalies	430
	23.7	Summary	433
Х	N S	Meta-Learning, Automated Learner Selection, Feature Selection, and Parameter Optimization	437
24	Usin Jovar 24.1 24.2	g RapidMiner for Research: Experimental Evaluation of Learners nović Miloš, Vukićević Milan, Delibašić Boris, and Suknović Milija Introduction	439 439 440
		24.2.1 Sources of Variation and Control	440

xv

	24.2.2 Example of an Experimental Setup	441
24.3	Experimental Evaluation in RapidMiner	442
	24.3.1 Setting Up the Evaluation Scheme	442
	24.3.2 Looping Through a Collection of Datasets	443
	24.3.3 Looping Through a Collection of Learning Algorithms	445
	24.3.4 Logging and Visualizing the Results	445
	24.3.5 Statistical Analysis of the Results	447
	24.3.6 Exception Handling and Parallelization	449
	24.3.7 Setup for Meta-Learning	450
24.4	Conclusions	452

Subject Index											•	•				•							455
Operator Index								•		•	•		•		•				•				463

xvi

Foreword

Case Studies Are for Communication and Collaboration

Data mining or data analysis in general has become more and more important, since large amounts of data are available and open up new opportunities for enhanced empirical sciences, planning and control, and targeted marketing and information services. Fortunately, theoretically well-based methods of data analysis and their algorithmic implementations are available. Experienced analysts put these programs to good use in a broad variety of applications. However, the successful application of algorithms is an art! There is no mapping from application tasks to algorithms, which could determine the appropriate chain of operators leading from the given data to the desired result of the analysis—but there are examples of such processes. Case studies are an easy way of communicating smart application design. This book is about such case studies in the field of data analysis.

Analysts are interested in the work of others and curiously inspect new solutions in order to stay up to date. Case studies are an excellent means to communicate best-practice compositions of methods to form a successful process. Cases are also well-suited to storing results for further use. A case is then used as a blueprint for further applications¹. This eases the development of an application to quite some extent.

Another good use of case studies is to ease the communication between application domain experts and data mining specialists. The case shows what could already be achieved and inspire future cases². This allows to frame new applications and to illustrate possible results. For those experts who want to set up a data mining project, it is a valuable justification.

Finally, teaching means communication. Teaching data mining is not complete without reference to case studies, either. Rapid-I offers at their website, http://rapid-i.com, video tutorials (webinars), white papers, manuals—a large variety of documentation with many illustrations. Offering case studies is now a further step into communicating not only the facilities of the system, but its use in real-world applications. The details of complex data analysis processes help those who want to become data analysts.

In summary, case studies support the collaboration of data analysts among themselves, the communication of data analysts and application experts, the interaction between experienced and beginners. Now, how can complex data mining processes be communicated, exchanged, and illustrated? An easy-to-understand view of the process is to abstract away the programming details. As is explained in the following, RapidMiner offers this.

¹T. Euler. Publishing operational models of data mining case studies. In B. Kitts, G. Melli, and K. Rexer, editors, *Procs. of the 1st Workshop on Data Mining Case Studies*, held at IEEE ICDM, pages 99-106, 2005.

²G. Melli, X Wu, P. Beinat, F. Bonchi, L. Cao, Rong Dan, C. Faloutsos, B. Kitts, B. Goethals, G. McLachlan, J. Pei, A. Srivastava, and O. Zaiane. Top10 data mining case studies. *Int. J. Information Technology and Decision Making*, 11(2):389-400, 2012.

RapidMiner

RapidMiner is a system which supports the design and documentation of an overall data mining process. It offers not only an almost comprehensive set of operators, but also structures that express the control flow of the process.

- Nesting operator chains were characteristic of RapidMiner (Yale) from the very beginning. This allows us to have a small number of operators at one level, each being expanded at the level below by simply clicking on the lower right area of the operator box.
- An example set can be multiplied for different processes that are executed in parallel and then be unified again. Sets of rows of different example sets of the same set of attributes can be appended. Hence, the example set that is used by some learning method can flexibly be modified.
- The cross validation is one of the most popular nested operators. The training set is split into n parts and, in a loop, n-1 parts are put to training and the remaining part to testing, so that the performance on the test set can be averaged over a range of different example sets from the same domain. The operator X-VALIDATION is used in most of the case studies in order to achieve sensible performance evaluations.
- Several loop operators can be specified for an application. The LOOP PARAMETERS operator repeatedly executes some other operators. The parameters of the inner operators as well as conditions controlling the loop itself tailor the operator to the desired control flow.
- Wrapper approaches wrap a performance-based selection around a learning operator. For instance, those feature sets are selected for which a learner's performance is best. The wrapper must implement some search strategy for iterating over sets and for each set call a learning algorithm and its evaluation in terms of some performance measure.

These structures are similar to notions of programming languages, but no programming is necessary – it is just drag, drop, and click! Visually, the structures are shown by boxes which are linked or nested. This presentation is easy to understand.

Only a small (though decisive) part of an overall data mining process is about model building. Evaluating and visualizing the results is the concluding part. The largest part is the pre-processing.

- It starts with reading in the data and declaring the meta data. RapidMiner supports many data formats and offers operators for assigning not only value domains of attributes (attribute type), but also their role in the learning process.
- The inspection of the data through diverse plots is an important step in developing the case at hand. In many case studies, this step is not recorded, since after the exploration it is no longer necessary. The understanding of the analysis task and the data leads to the successful process.
- Operators that change the given representation are important to bridge the gap between the given data and the input format that a learner needs. Most analysts have a favorite learning method and tweak the given data until they suit this algorithm well. If frequent set mining is the favorite, the analyst will transform m nominal values of

xviii

one attribute into m binomial attributes so that frequent set mining can be applied. If the attribute type requirements of a learner are not yet fulfilled, RapidMiner proposes fixes.

- The discretization of real-valued attributes into nominal- or binomial-valued attributes is more complex and, hence, RapidMiner offers a variety of operators for this task.
- Beyond type requirements, features extraction and construction allow learners to find interesting information in data which otherwise would be hidden. A very large collection of operators offers the transformation of representations. The text processing plug-in, the value series plug-in, and the image processing plug-in are specifically made for the pre-processing of texts, time series or value series in general, and images.
- The feature selection plug-in automatically applies user-specified criteria to design the best feature set for a learning task. Moreover, it evaluates the selected features with respect to stability. For real-world applications, it is important that good performance is achieved at any sample of data. It is not sufficient that the selected features allow a good performance on average in the cross-validation runs, but it must be guaranteed that the features allow a sufficiently good performance on every data sample.

Given the long operator chains and nested processes in data mining, the aspect of documentation becomes indispensable. The chosen parameters of, e.g., discretization, the particular feature transformations, and the criteria of feature selection are stored with the Rapid-Miner process. The metadata characterize the data at any state in the process. Hence, its result is explainable and reproducible.

In this book, case studies communicate how to analyze databases, text collections, and image data. The favorite learning tasks are classification and regression with the favorite learning method being support vector machines followed by decision trees. How the given data are transformed to meet the requirements of the method is illustrated by pictures of RapidMiner. The RapidMiner processes and datasets described in the case studies are published on the companion web page of this book. The inspiring applications may be used as a blueprint and a justification of future applications.

Prof. Dr. Katharina Morik (TU Dortmund, Germany)

Preface

Data and the ability to make the best use of it are becoming more and more crucial for today's and tomorrow's companies, organizations, governments, scientists, and societies to tackle everyday challenges as well as complex problems and to stay competitive. Data mining, predictive analytics, and business analytics leverage these data, provide unprecedented insights, enable better-informed decisions, deliver forecasts, and help to solve increasingly complex problems. Companies and organizations collect growing amounts of data from all kinds of internal and external sources and become more and more data-driven. Powerful tools for mastering data analytics and the know-how to use them are essential to not fall behind, but to gain competitive advantages, and to increase insights, effectiveness, efficiency, growth, and profitability.

This book provides an introduction to data mining and business analytics, to the most powerful and flexible open source software solutions for data mining and business analytics, namely *RapidMiner and RapidAnalytics*, and to many application *use cases* in scientific research, medicine, industry, commerce, and diverse other sectors. RapidMiner and RapidAnalytics and their extensions used in this book are all freely available as open source software community editions and can be downloaded from http://www.RapidMiner.com

Each chapter of this book describes an application *use case*, how to approach it with *data mining methods*, and how to implement it with RapidMiner and RapidAnalytics. These application-oriented chapters do not only provide you with the necessary analytics knowhow to solve these problems and tasks, but also with easy-to-follow reproducible step-by-step descriptions for accomplishing this with RapidMiner and RapidAnalytics. The datasets and RapidMiner processes used in this book are available from the companion web page of this book:

http://www.RapidMinerBook.com

This application-oriented analytics use case collection will quickly enable you to solve similar problems effectively yourself. The case studies can serve as blue prints for your own data mining applications.

What Is Data Mining? What Is It Good for, What Are Its Applications, and What Does It Enable Me to Do?

While technology enables us to capture and store ever larger quantities of *data*, finding relevant *information* like underlying patterns, trends, anomalies, and outliers in the data and summarizing them with simple understandable and robust quantitative and qualitative models is a grand challenge. Data mining helps to discover underlying structures in the data, to turn *data* into *information*, and information into *knowledge*. Emerged from mathematics,

statistics, logic, computer science, and information theory, data mining and machine learning and statistical learning theory now provide a solid theoretical foundation and powerful methods to master this challenge. Data mining is the extraction of implicit, previously unknown, and potentially useful information from data. The automatically extracted models provide insight into customer behavior and into processes generating data, but can also be applied to, for example, automatically classify objects or documents or images into given categories, to estimate numerical target variables, to predict future values of observed time series data, to recommend products, to prevent customer churn, to optimize direct marketing campaigns, to forecast and reduce credit risk, to predict and prevent machine failures before they occur, to automatically route e-mail messages based on their content and to automatically detect e-mail spam, and to many other tasks where data helps to make better decisions or even to automate decisions and processes. Data mining can be applied not only to structured data from files and databases, but text mining extends the applicability of these techniques to unstructured data like texts from documents, news, customer feedback, e-mails, web pages, Internet discussion groups, and social media. Image mining, audio mining, and video mining apply these techniques to even further types of data.

Why Should I read This Book? Why Case Studies? What Will I learn? What Will I Be Able to Achieve?

This book introduces the most important machine learning algorithms and data mining techniques and enables you to use them in real-world applications. The open source software solutions RapidMiner and RapidAnalytics provide implementations for all of these algorithms and a powerful and flexible framework for their application to all kinds analytics tasks. The book and these software tools cover all relevant steps of the data mining process from data loading, transformation, integration, aggregation, and visualization via modeling, model validation, performance evaluation, model application and deployment, to automated feature selection, automated parameter and process optimization, and integration with other tools like, for example, the open source statistics package R or into your IT infrastructure via web services. The book and the tools also extensively cover the analysis of unstructured data including text mining and image mining.

The application-oriented focus of this book and the included use cases provide you with the know-how and blueprints to quickly learn how to efficiently design data mining processes and how to successfully apply them to real-world tasks. The book not only introduces you to important machine learning methods for tasks like clustering, classification, regression, association and recommendation rule generation, outlier and anomaly detection, but also to the data preprocessing and transformation techniques, which often are at least as crucial for success in real-world applications with customer data, product data, sales data, transactional data, medical data, chemical molecule structure data, textual data, web page data, image data, etc. The use cases in this book cover domains like retail, banking, marketing, communication, education, security, medicine, physics, and chemistry and tasks like direct marketing campaign optimization, product affinity scoring, customer churn prediction and prevention, automated product recommendation, increasing sales volume and profits by cross-selling, automated video lecture recommendation, intrusion detection, fraud detection, credit approval, automated text classification, e-mail and mobile phone text message spam detection, automated language identification, customer feedback and hotel review analysis, image classification, image feature extraction, automated feature selection, clustering students in higher education and automated study program recommendation, ranking school applicants, teaching assistant evaluation, pharmaceutical molecule activity prediction, medical research, biochemical research, neutrino physics research, and data mining research.

What Are the Advantages of the Open Source Solutions RapidMiner and RapidAnalytics Used in This Book?

RapidMiner and RapidAnalytics provide an integrated environment for all steps of the data mining process, an easy-to-use graphical user interface (GUI) for the interactive data mining process design, data and results visualization, validation and optimization of these processes, and for their automated deployment and possible integration into more complex systems. RapidMiner enables one to design data mining processes by simple drag and drop of boxes representing functional modules called operators into the process, to define data flows by simply connecting these boxes, to define even complex and nested control flows, and all without programming. While you can seamlessly integrate, for example, R scripts or Groovy scripts into these processes, you do not need to write any scripts, if you do not want to. RapidMiner stores the data mining processes in a machine-readable XML format, which is directly executable in RapidMiner with the click of a button, and which along with the graphical visualization of the data mining process, makes it easy to execute, to validate, to automatically optimize, to reproduce, and to automate.

Their broad functionality for all steps of the data mining process and their flexibility make RapidMiner and RapidAnalytics the tools of choice. They optimally support all steps of the overall data mining process and the flexible deployment of the processes and results within their framework and also integrated into other solutions via web services, Java API, or command-line calls. The process view of data mining eases the application to complex real-world tasks and the structuring and automation of even complex highly nested data mining processes. The processes also serve as documentation and for the reproducibility and reusability of scientific results as well as business applications. The open source nature of RapidMiner and RapidAnalytics, their numerous import and export and web service interfaces, and their openness, flexibility, and extendibility by custom extensions, operators, and scripts make them the ideal solutions for scientific, industrial, and business applications. Being able to reproduce earlier results, to reuse previous processes, to modify and adapt them or to extend them with customized or self-developed extensions means a high value for research, educational training, and industrial and business applications, RapidMiner allows you to quickly build working prototypes and also quickly deploy them on real data of all types including files, databases, time series data, text data, web pages, social media, image data, audio data, web services, and many other data sources.

What Is the Structure of This Book and Which Chapters Should I Read?

The first chapter of this book introduces the basic concepts of data mining and machine learning, common terms used in the field and throughout this book, and the decision tree modeling technique as a machine learning technique for classification tasks. The second chapter gives you an introductory tour through the RapidMiner graphical user interface (GUI) and how to use it to define data mining processes. In case you are already familiar with data mining and RapidMiner, you can skip these two chapters. However, if you are a novice in the field or regarding the software, these first two chapters are highly recommended and will give you a quick start in both data mining and RapidMiner. All following chapters provide a use case each and introduce additional data mining concepts and RapidMiner operators needed to solve the task at hand.

The Chapters 3 to 6 describe classification use cases and introduce the k-nearest neighbors (k-NN) and Naïve Bayes learning algorithms. Chapter 3 applies k-NN for the evaluation of teaching assistants. In Chapter 4 k-NN is used to classify different glass types based on chemical components and the RapidMiner process is extended by Principal Component Analysis (PCA) to better preprocess the data and to improve the classification accuracy. Chapter 5 explains Naïve Bayes as an algorithm for generating classification models and uses this modeling technique to generate a credit approval model to decide whether a credit loan for which a potential or existing customer applies should be approved or not, i.e. whether it is likely that the customer will pay back the credit loan as desired or not. Chapter 6 uses Naïve Bayes to rank applications for nursery schools, introduces the RapidMiner operator for importing Excel sheets, and provides further explanations of Naïve Bayes.

Chapter 7 addresses the task of product affinity-based marketing and optimizing a direct marketing campaign. A bank has introduced a new financial product, a new type of current (checking) account, and some of its customers have already opened accounts of the new type, but many others have not done so yet. The bank's marketing department wants to push sales of the new account by sending direct mail to customers who have not yet opted for it. However, in order not to waste efforts on customers who are unlikely to buy, they would like to address only those customers with the highest affinity for the new product. Binary classification is used to predict for each customer, whether they will buy the product, along with a confidence value indicating how likely each of them is to buy the new product. Customers are then ranked by this confidence value and the 20% with the highest expected probability to buy the product are chosen for the campaign.

Following the *CRoss-Industry Standard Process for Data Mining (CRISP-DM)* covering all steps from business understanding and data understanding via data preprocessing and modeling to performance evaluation and deployment, this chapter first describes the task, the available data, how to extract characteristic customer properties from the customer data, their products and accounts data and their transactions, which data preprocessing to apply to balance classes and aggregate information from a customer's accounts and transactions into attributes for comparing customers, modeling with binary classification, evaluating the predictive accuracy of the model, visualizing the performance of the model using Lift charts and ROC charts, and finally ranking customers by the predictive accuracy of several learning algorithms including Decision Trees, Linear Regression, and Logistic Regression is compared and visualized comparing their ROC charts. Automated attribute weight and parameter optimizations are deployed to maximize the prediction accuracy and thereby the customer response, sales volume, and profitability of the campaign. Similar processes campaing be used for *customer churn prediction* and addressing the customers predicted to churn in a campaign with special offers trying to prevent them from churning.

Chapters 8 to 10 describe three different approaches to building recommender systems. Product recommendations in online shops like Amazon increase the sales volume per customer by cross-selling, i.e., by selling more products per customer by recommending products that the customer may also like and buy.

The recommendations can be based on product combinations frequently observed in market baskets in the past. Products that co-occurred in many purchases in the past are assumed to be also bought together frequently in the future. *Chapter 8* describes how to generate such association rules for product recommendations from shopping cart data using the FP-Growth algorithm. Along the way, this chapter also explains how to import product sales data from CSV files and from retailers' databases and how to handle data quality issues and missing values.

Chapter 9 introduces the RapidMiner Extension for Recommender Systems. This extension allows building more sophisticated recommendation systems than described in the previous chapter. The application task in this chapter is to recommend appropriate video lectures to potential viewers. The recommendations can be based on the content of the lectures or on the viewing behavior or on both. The corresponding approaches are called content-based, collaborative, and hybrid recommendation, respectively. Content-based recommendations can be based on attributes or similarity and collaborative recommendation systems deploy neighborhoods or factorization. This chapter explains, evaluates, and compares these approaches. It also demonstrates how to make RapidMiner processes available as RapidAnalytics web services, i.e., how to build a recommendation engine and make it available for real-time recommendations and easy integration into web sites, online shops, and other systems via web services.

A third way of building recommender systems in RapidMiner is shown in *Chapter 10*, where classification algorithms are used to recommend the best-fitting study program for higher-education students based on their predicted success for different study programs at a particular department of a particular university. The idea is an early analysis of students' success on each study program and the recommendation of a study program where a student will likely succeed. At this university department, the first year of study is common for all students. In the second year, the students select their preferred study program among several available programs. The attributes captured for each graduate student describe their success in the first-year exams, their number of points in the entrance examination, their sex, and their region of origin. The target variable is the average grade of the student at graduation, which is discretized into several categories. The prediction accuracy of several classification learning algorithms, including Naïve Bayes, Decision Trees, Linear Model Tree (LMT), and CART (Classification and Regression Trees), is compared for the prediction of the student's success as measured by the discretized average grade. For each student, the expected success classes for each study program is predicted and the study program with the highest predicted success class is recommended to the student. An optimization loop is used to determine the best learning algorithm and automated feature selection is used to find the best set of attributes for the most accurate prediction. The RapidMiner processes seamlessly integrate and compare learning techniques implemented in RapidMiner with learning techniques implemented in the open source data mining library Weka, thanks to the Weka extension for RapidMiner that seamlessly integrates all Weka learners into RapidMiner.

Chapter 11 provides an introduction to clustering, to the k-Means clustering algorithm, to several cluster validity measures, and to their visualizations. Clustering algorithms group cases into groups of similar cases. While for classification, a training set with examples with predefined categories is necessary for training a classifier to automatically classify new cases

into one of the predefined categories, clustering algorithms need no labeled training examples with predefined categories, but automatically group unlabeled examples into clusters of similar cases. While the predictive accuracy of classification algorithms can be easily measured by comparing known category labels of known examples to the categories predicted by the algorithm, there are no labels known in advance in the case of clustering. Hence it is more difficult to achieve an objective evaluation of a clustering result. Visualizing cluster validity measures can help humans to evaluate the quality of a set of clusters. This chapter uses k-Means clustering on a *medical dataset to find groups of similar E-Coli bacteria with regards to where protein localization occurs in them* and explains how to judge the quality of the clusters found using visualized cluster validity metrics. Cluster validity measures implemented in the open source statistics package R are seamlessly integrated and used within RapidMiner processes, thanks to the *R extension for RapidMiner*.

Chapter 12 applies clustering to automatically group higher education students. The dataset corresponds to the one already described in Chapter 10, but now the task is to find groups of similarly performing students, which is achieved with automated clustering techniques. The attributes describing the students may have missing values and different scales. Hence data preprocessing steps are used to replace missing values and to normalize the attribute values to identical value ranges. A parameter loop automatically selects and evaluates the performance of several clustering techniques including k-Means, k-Medoids, Support Vector Clustering (SVC), and DBSCAN.

Chapters 13 to 15 are about text mining applications. Chapter 13 gives an introduction to text mining, i.e., the application of data mining techniques like classification to text documents like e-mail messages, mobile phone text messages (SMS = Short Message Service) or web pages collected from the World-Wide Web. In order to detect text message spam, preprocessing steps using the RapidMiner text processing extension transform the unstructured texts into document vectors of equal length, which make the data applicable to standard classification techniques like Naïve Bayes, which is then trained to automatically separate legitimate mobile phone text messages from spam messages.

The second text mining use case uses classification to automatically identify the language of a text based on its characters, character sequences, and/or words. Chapter 14 discusses character encodings of different European, Arabic, and Asian languages. The chapter describes different text representations by characters, by tokens like words, and by character sequences of a certain length also called n-grams. The transformation of document texts into document vectors also involves the weighting of the attributes by term frequency and document frequency-based metrics like TF/IDF, which is also described here. The classification techniques Naïve Bayes and Support Vector Machines (SVM) are then trained and evaluated on four different multi-lingual text corpora including for example dictionary texts from Wikipedia and book texts from the Gutenberg project. Finally, the chapter shows how to make the RapidMiner language detection available as web service for the automated language identification of web pages via RapidAnalytics web services.

Chapter 15 analyses hotel review texts and ratings by customers collected from the TripAdvisor web page. Frequently co-occurring words in the review texts are found using FP-Growth and association rule generation and visualized in a word-association graph. In a second analysis, the review texts are clustered with k-Means, which reveals groups of similar texts. Both approaches provide insights about the hotels and their customers, i.e., about topics of interest and of complaints, quality and service issues, likes, dislikes, and preferences, and could similarly be applied to all kinds of textual reviews and customer feedback.

Chapter 16 describes a data mining use case in *astroparticle physics*, the application of automated *classification* and automated *feature selection* in *neutrino astronomy* to separate a small number of neutrinos from a large number of background noise particles or signals

xxvi

(muons). One of the main scientific goals of neutrino telescopes is the detection of neutrinos originating from astrophysical sources as well as a precise measurement of the energy spectrum of neutrinos produced in cosmic ray air showers in the Earth's atmosphere. These so-called atmospheric neutrinos, however, are hidden in a noisy background of atmospheric muons produced in air showers as well. The first task in rejecting this background is the selection of upward-going tracks since the Earth is opaque to muons but can be traversed by neutrinos up to very high energies. This procedure reduces the background by roughly three orders of magnitude. For a detailed analysis of atmospheric neutrinos, however, a very clean sample with purity larger than 95% is required. The main source of remaining background at this stage are muon tracks, falsely reconstructed as upward going. These falsely reconstructed muon tracks still dominate the signal by three orders of magnitude and have to be rejected by the use of straight cuts or multivariate methods. Due to the ratio of noise (muons) and signal (neutrinos), about 10,000 particles need to be recorded in order to catch about 10 neutrinos. Hence, the amount of data delivered by these experiments is enormous and it must be processed and analyzed within a proper amount of time. Moreover, data in these experiments are delivered in a format that contains more than 2000 attributes originating from various reconstruction algorithms. Most of these attributes have been reconstructed from only a few physical quantities. The direction of a neutrino event penetrating the detector at a certain angle can, for example, be reconstructed from a pattern of light that is initiated by particles produced by an interaction of the neutrino close to or even in the detector. Due to the fact that all of the 2000 reconstructed attributes are not equally well suited for classification, the first task in applying data mining techniques in neutrino astronomy lies in finding a good and reliable representation of the dataset in fewer dimensions. This is a task which very often determines the quality of the overall data analysis. The second task is the training and evaluation of a stable learning algorithm with a very high performance in order to separate signal and background events. Here, the challenge lies in the biased distribution of many more background noise (negative) examples than there are signals (positive) examples. Handling such skewed distributions is necessary in many real-world problems. The application of RapidMiner in neutrino astronomy models the separation of neutrinos from background as a two-step process, accordingly. In this chapter, the *feature or attribute selection* is explained in the first part and the *training of* selecting relevant events from the masses of incoming data is explained in the second part. For the feature selection, the Feature Selection Extension for RapidMiner is used and a wrapper cross-validation to evaluate the performance of the feature selection methods. For the selection of the relevant events, *Random Forests* are used as classification learner.

Chapter 17 provides an introduction to *medical data mining*, an overview of methods often used for classification, regression, clustering, and association rules generation in this domain, and two application use cases with data about patients suffering from *carpal tunnel syndrome* and *diabetes*, respectively.

In the study of the *carpal tunnel syndrome (CTS)*, thermographic images of hands were collected for constructing a predictive classification model for CTS, which could be helpful when looking for a non-invasive diagnostic method. The temperatures of different areas of a patient's hand were extracted from the image and saved in the dataset. Using a RapidMiner preprocessing operator for aggregation, the temperatures were averaged for all segments of the thermal images. Different machine learning algorithms including *Artificial Neural Network* and *Support Vector Machines (SVM)* were evaluated for generating a classification model capable of diagnosing CTS on the basis of very discrete temperature differences that are invisible to the human eye in a thermographic image.

In the study of *diabetes*, various research questions were posed to evaluate the level of knowledge and overall perceptions of diabetes mellitus type 2 (DM) within the older population in North-East Slovenia. As a chronic disease, diabetes represents a substantial

xxviii

burden for the patient. In order to accomplish good self-care, patients need to be qualified and able to accept decisions about managing the disease on a daily basis. Therefore, a high level of knowledge about the disease is necessary for the patient to act as a partner in managing the disease. Various research questions were posed to determine what the general knowledge about diabetes is among diabetic patients 65 years and older, and what the difference in knowledge about diabetes is with regard to the education and place of living on (1) diet, (2) HbA1c, (3) hypoglycemia management, (4) activity, (5) effect of illness and infection on blood sugar levels, and (6) foot care. A hypothesis about the level of general knowledge of diabetes in older populations living in urban and rural areas was predicted and verified through the study. A cross-sectional study of older (age >65 years), non-insulin dependent patients with diabetes mellitus type 2 who visited a family physician, DM outpatient clinic, a private specialist practice, or were living in a nursing home was implemented. The Slovenian version of the Michigan Diabetes Knowledge test was then used for data collection. In the data preprocessing, missing values in the data were replaced, before k-means clustering was used to find groups of similar patients, for which then a decision tree learner was used to find attributes discriminating the clusters and generate a classification model for the clusters. A grouped ANOVA (ANalysis Of VAriances) statistical test verified the hypothesis that there are differences in the level of knowledge about diabetes in rural populations and city populations in the age group of 65 years and older.

Chapter 18 covers a use case relevant in chemistry and the pharmaceutical industry. The RapidMiner Extension PaDEL (Pharmaceutical Data Exploration Laboratory) developed at the University of Singapore is deployed to calculate a variety of molecular properties from the 2-D or 3-D molecular structures of chemical compounds. Based on these molecular property vectors, RapidMiner can then generate predictive models for predicting chemical, biochemical, or biological properties based on molecular properties, which is a frequently encountered task in theoretical chemistry and the pharmaceutical industry. The combination of RapidMiner and PaDEL provides an open source solution to generate prediction systems for a broad range of biological properties and effects.

One application example in *drug design* is the *prediction of effects and side effects* of a new drug candidate before even producing it, which can help to avoid testing many drug candidates that probably are not helpful or possibly even harmful and thereby help to focus research resources on more promising drug candidates. With PaDEL and RapidMiner, properties can be calculated for any molecular structure, even if the compound is not physically accessible. Since both tools are open source and can compute the properties of a molecular structure quickly, this allows significant reduction in cost and an increase in speed of the development of new chemical compounds and drugs with desired properties, because more candidate molecules can be considered automatically and fewer of them need to be actually generated and physically, chemically, or biologically tested.

The combination of data mining (RapidMiner) and a tool to handle molecules (PaDEL) provides a convenient and user-friendly way to generate accurate relationships between chemical structures and any property that is supposed to be predicted, mostly biological activities. Relationships can be formulated as qualitative structure-property relationships (SPRs), qualitative structure-activity relationships (SARs) or quantitative structure-activity relationships (SARs) or quantitative structure-activity relationships (SARs). SPR models aim to highlight associations between molecular structures and a target property, such as lipophilicity. SAR models correlate an activity with structural properties and QSAR models quantitatively predict an activity. Models are typically developed to predict properties that are difficult to obtain, impossible to measure, require time-consuming experiments, or are based on a variety of other complex properties. They may also be useful to predict complicated properties using several simple properties. The PaDEL extension enables RapidMiner to directly read and handle molecular structures, calculate their molecular properties, and to then correlate them to and generate predictive

models for chemical, biochemical, or biological properties of these molecular structures. In this chapter *linear regression* is used as a QSAR modeling technique to predict chemical properties with RapidMiner based on molecular properties computed by PaDEL.

Chapter 19 describes a second Quantitative Structure-Activity Relationship (QSAR) use case relevant in chemistry and the pharmaceutical industry, the identification of novel functional inhibitors of acid sphingomyelinase (ASM). The use case in this chapter is based on the previous chapter and hence you should first read Chapter 18 before reading this chapter. In the data preprocessing step, the PaDEL (Pharmaceutical Data Exploration Laboratory) extension for RapidMiner described in the previous chapter is again used to compute molecular properties from given molecular 2-D or 3-D structures. These properties are then used to predict ASM inhibition. Automated feature selection with backward elimination is used to reduce the number of properties to a relevant set for the prediction task, for which a classification learner, namely Random Forests, generates the predictive model that captures the structure- and property-activity relationships.

The process of drug design from the biological target to the drug candidate and, subsequently, the approved drug has become increasingly expensive. Therefore, strategies and tools that reduce costs have been investigated to improve the effectiveness of drug design. Among them, the most time-consuming and cost-intensive steps are the selection, synthesis, and experimental testing of the drug candidates. Therefore, numerous attempts have been made to reduce the number of potential drug candidates for experimental testing. Several methods that rank compounds with respect to their likelihood to act as an active drug have been developed and applied with variable success. In silico methods that support the drug design process by reducing the number of promising drug candidates are collectively known as virtual screening methods. Their common goal is to reduce the number of drug candidates subjected to biological testing and to thereby increase the efficacy of the drug design process.

This chapter demonstrates an in silico method to predict biological activity based on RapidMiner data mining work flows. This chapter is based on the type of chemoinformatic predictions described in the previous chapter based on chemoinformatic descriptors computed by PaDEL. Random Forests are used as a predictive model for predicting the molecular activity of a molecule of a given structure, for which PaDEL is used to compute molecular structural properties, which are first reduced to a smaller set by automated *attribute weighting* and *selecting the attributes* with the highest weights according to several weighting criteria and which are reduced to an even smaller set of attributes by automated *attribute selection* using a *Backward Elimination* wrapper. Starting with a large number of properties for the example set, a *feature selection* vastly reduces the number of attributes before the systematic backward elimination search finds the most predictive model for the feature generation. Finally, a validation is performed to avoid over-fitting and the benefits of *Y-randomization* are shown.

Chapter 20 introduces the RapidMiner IMage Mining (IMMI) Extension and presents some introductory image processing and image mining use cases. Chapter 21 provides more advanced image mining applications.

Given a set of images in a file folder, the image processing task in the first use case in *Chapter 20* is to adjust the contrast in all images in the given folder and to store the transformed images in another folder. The IMMI extension provides RapidMiner operators for *reading and writing images*, which can be used within a RapidMiner loop iterating over all files in the given directory, *adjusting the contrast* of each of these images, for example, using a histogram equalization method. Then the chapter describes image *conversions between color and gray-scale images* and different *feature extraction* methods, which convert image data in unstructured form into a tabular form. Feature extraction algorithms for images can be divided into three basic categories: local-level, segment-level, and global-level feature extraction.

The term *local-level* denotes that information is mined from given points (locations) in the image. *Local-level feature extraction* is suitable for *segmentation*, *object detection* or *area detection*. From each point in the image, it is possible to extract information like pixel gray value, minimal or maximal gray value in a specified radius, value after applying kernel function (blurring, edge enhancements). Examples of utilization of such data are the trainable segmentation of an image, point of interest detection, and object detection.

The term *segment-level* denotes feature extraction from segments. Many different segmentation algorithms exist, such as k-means, watershed, or statistical region merging. Segment level feature extraction algorithms extract information from the whole segments. Examples of such features are mean, median, lowest and highest gray value, circularity, and eccentricity. In contrast to local-level features, it does not take into consideration only a single point and its neighborhood, however, it considers the whole segment and its properties like shape, size, and roundness. With the use of knowledge about the size or shape of target objects, it is for example possible to *select or remove objects according to their size or shape*.

The global-level denotes feature extraction from the whole image, for example, mean color, dominant color, maximal gray value, minimal gray value, variance of pixels, number of edges etc. Unlike the local or segment level, the global level segmentation is not suitable for points or areas identification or segmentation. Rather, it is suitable for classification of images and determining properties of the image as a whole.

Chapter 20 provides examples demonstrating the use of local-level, segment-level, and global-level feature extraction. Local-level feature extraction is used for trainable image segmentation with radial-basis function (RBF) Support Vector Machines (SVM). Segment-level feature extraction and trainable segment selection reveal interesting segment properties like size and shape for image analysis. With the help of global-level feature extraction, images are classified into pre-defined classes. In the presented use case, two classes of images are distinguished automatically: images containing birds and images containing sunsets. To achieve this, global features like dominant color, minimal intensity, maximal intensity, percent of edges, etc. are extracted and based on those, an image classifier is trained.

Chapter 21 presents advanced image mining applications using the RapidMiner IMage Mining (IMMI) Extension introduced in the previous chapter. This chapter demonstrates several examples of the use of the IMMI extension for image processing, image segmentation, feature extraction, pattern detection, and image classification. The first application extracts global features from multiple images to enable automated image classification. The second application demonstrates the Viola-Jones algorithm for pattern detection. And the third process illustrates the image segmentation and mask processing.

The classification of an image is used to identify which group of images a particular image belongs to. An automated image classifier could, for example, be used to distinguish different scene types like nature versus urban environment, exterior versus interior, images with and without people, etc. *Global features* are usually used for this purpose. These features are calculated from the whole image. The key to a correct classification is to find the features that differentiate one class from other classes. Such a feature can be, for example, the dominant color in the image. These features can be calculated from the original image or from an image after pre-processing like Gaussian blur or edge detection.

Pattern detection searches known patterns in images in the images, where approximate fits of the patterns may be sufficient. A good algorithm for detection should not be sensitive to the size of the pattern in the image or its position or rotation. One possible approach is to use a histogram. This approach compares the histogram of the pattern with the histogram of a selected area in the image. In this way, the algorithm passes step by step through

XXX

the whole image, and if the match of histograms is larger than a certain threshold, the area is declared to be the sought pattern. Another algorithm, which is described in this chapter, is the Viola-Jones algorithm. The classifier is trained with positive and negative image examples. Appropriate features are selected using the AdaBoost algorithm. An image is iterated during pattern detection using a window with increasing size. Positive detections are then marked with a square area of the same size as the window. The provided example application uses this process to detect the cross-sectional artery in an ultrasound image. After detection, the images can be used to measure the patient's pulse if taken from a video or stream of time-stamped images.

The third example application demonstrates *image segmentation* and *feature extraction*: Image segmentation is often used for the *detection of different objects* in the image. Its task is to split the image into parts so that the individual segments correspond to objects in the image. In this example, the identified segments are combined with *masks* to remove the background and focus on the object found.

Chapter 22 introduces the RapidMiner Extension for Instance Selection and Prototypebased Rule (ISPR) induction. It describes the instance selection and prototype construction methods implemented in this extension and applies them to accelerate 1-NN classification on large datasets and to perform *outlier elimination* and *noise reduction*. The datasets analyzed in this chapter include several *medical datasets* for classifying patients with respect to certain medical conditions, i.e., diabetes, heart diseases, and breast cancer, as well as an e-mail spam detection dataset. The chapter describes a variety of prototype selection algorithms including k- Nearest-Neighbors (k-NN), Monte-Carlo (MC) algorithm, Random Mutation Hill Climbing (RMHC) algorithm, Condensed Nearest-Neighbor (CNN), Edited Nearest-Neighbor (ENN), Repeated ENN (RENN), Gabriel Editing proximity graph-based algorithm (GE selection), Relative Neighbor Graph algorithm (RNG selection), Instance-Based Learning (IBL) algorithm (IB3 selection), Encoding Length Heuristic (ELH selection), and combinations of them and compares their performance on the datasets mentioned above. Prototype construction methods include all algorithms that produce a set of instances at the output. The family contains all prototype-based clustering methods like k-Means, Fuzzy C-Means (FCM), and Vector Quantization (VQ) as well as the Learning Vector Quantization (LVQ) set of algorithms. The price for the speed-up of 1-NN by instance selection is visualized by the drop in predictive accuracy with decreasing sample size.

Chapter 23 gives an overview of a large range of anomaly detection methods and introduces the RapidMiner Anomaly Detection Extension. Anomaly detection is the process of finding patterns in a given dataset which deviate from the characteristics of the majority. These outstanding patterns are also known as anomalies, outliers, intrusions, exceptions, misuses, or fraud. Anomaly detection identifies single records in datasets which significantly deviate from the normal data. Application domains among others include *network security*, intrusion detection, computer virus detection, fraud detection, misuse detection, complex system supervision, and finding suspicious records in medical data. Anomaly detection for fraud detection is used to *detect fraudulent credit card transactions* caused by stolen credit cards, fraud in Internet payments, and suspicious transactions in financial accounting data. In the medical domain, anomaly detection is also used, for example, for detecting tumors in medical images or monitoring patient data (electrocardiogram) to get early warnings in case of life-threatening situations. Furthermore, a variety of other specific applications exists such as anomaly detection in surveillance camera data, fault detection in complex systems or detecting forgeries in the document forensics. Despite the differences of the various application domains, the basic principle remains the same. Multivariate normal data needs to be modeled and the few deviations need to be detected, preferably with a score indicating their "outlierness", i.e., a score indicating their extent of being an outlier. In case xxxii

of a univariate data, such an *outlier factor* could for example be the number of standard deviations by which an outlier differs from the mean of this variable.

The overview of anomaly detection method provided in this chapter distinguishes three different types of anomalies, namely (1) point anomalies, which are single data records deviating from others, (2) contextual anomalies, which occur with respect to their context only, for example, with respect to time, and (3) collective anomalies, where a bunch of data points causes the anomaly. Most anomaly detection algorithms detect point anomalies only, which leads to the requirement of transforming contextual and collective anomalies to point anomaly problems using an appropriate pre-processing and thus generating processable data views. Furthermore, anomaly detection algorithms can be categorized with respect to their operation mode, namely (1) supervised algorithms with training and test data as used in traditional machine learning, (2) semi-supervised algorithms with the need of anomaly-free training data for one-class learning, and (3) unsupervised approaches without the requirement of any labeled data. Anomaly detection is, in most cases, associated with an unsupervised setup, which is also the focus of this chapter. In this context, all available unsupervised algorithms from the RapidMiner anomaly detection extension are described and the most well-known algorithm, the Local Outlier Factor (LOF) is explained in detail in order to get a deeper understanding of the approaches themselves. The unsupervised anomaly detection algorithms covered in this chapter include Grubbs' outlier test and noise removal procedure, k-NN Global Anomaly Score, Local Outlier Factor (LOF), Connectivity-Based Outlier Factor (COF), Influenced Outlierness (INFLO), Local Outlier Probability (LoOP), Local Correlation Integral (LOCI) and aLOCI, Cluster-Based Local Outlier Factor (CBLOF), and Local Density Cluster-Based Outlier Factor (LDCOF). The semi-supervised anomaly detection algorithms covered in this chapter include a one-class Support Vector Machine (SVM) and a two-step approach with clustering and distance computations for detecting anomalies.

Besides a simple example consisting of a two-dimensional mixture of Gaussians, which is ideal for first experiments, two real-world datasets are analyzed. For the unsupervised anomaly detection the player statistics of the NBA, i.e., a dataset with the NBA regularseason basketball player statistics from 1946 to 2009, are analyzed for outstanding players, including all necessary pre-processing. The UCI NASA shuttle dataset is used for illustrating how semi-supervised anomaly detection can be performed in RapidMiner to find suspicious states during a NASA shuttle mission. In this context, a Groovy script is implemented for a simple semi-supervised cluster-distance-based anomaly detection approach, showing how to easily extend RapidMiner by your own operators or scripts.

Chapter 24 features a complex data mining research use case, the performance evaluation and comparison of several classification learning algorithms including Naïve Bayes, k-NN, Decision Trees, Random Forests, and Support Vector Machines (SVM) across many different datasets. Nested process control structures for loops over datasets, loops over different learning algorithms, and cross validation allow an automated validation and the selection of the best model for each application dataset. Statistical tests like t-test and ANOVA test (ANalysis Of VAriance) determine whether performance differences between different learning techniques are statistically significant or whether they may be simply due to chance. Using a custom-built Groovy script within RapidMiner, meta-attributes about the datasets are extracted, which can then be used for *meta-learning*, i.e., for learning to predict the performance of each learner from a given set of learners for a given new dataset, which then allows the selection of the learner with the best expected accuracy for the given dataset. The performance of fast learners called landmarkers on a given new dataset and the metadata extracted from the dataset can be used for meta-learning to predict the performance of another learner on this dataset. The RapidMiner Extension for Pattern Recognition Engineering (PaREn) and its Automatic System Construction Wizard perform this kind of meta-learning for automated learner selection and a parameter optimization for a given dataset.

The *index at the end of the book* helps you to find explanations of data mining concepts and terms you would like to learn more about, use case applications you may be interested in, or reference use cases for certain modeling techniques or RapidMiner operators you are looking for. The *companion web page for this book* provides the RapidMiner processes and datasets deployed in the use cases:

http://www.RapidMinerBook.com

xxxiii

About the Editors

Markus Hofmann

Dr. Markus Hofmann is currently a lecturer at the Institute of Technology Blanchardstown in Ireland where he focuses on the areas of data mining, text mining, data exploration and visualisation as well as business intelligence. He holds a PhD degree from Trinity College Dublin, an MSc in Computing (Information Technology for Strategic Management) from the Dublin Institute of Technology and a BA in Information Management Systems. He has taught extensively at undergraduate and postgraduate level in the fields of Data Mining, Information Retrieval, Text/Web Mining, Data Mining Applications, Data Pre-processing and Exploration and Databases. Dr. Hofmann published widely at national as well as international level and specialised in recent years in the areas of Data Mining, learning object creation, and virtual learning environments. Further he has strong connections to the Business Intelligence and Data Mining sector both on an academic as well as industry level. Dr. Hofmann has worked as technology expert together with 20 different organisations in recent years including companies such as Intel. Most of his involvement was on the innovation side of technology services and products where his contributions had significant impact on the success of such projects. He is a member of the Register of Expert Panellists of the Irish Higher Education and Training Awards council, external examiner to two other third level institutes and a specialist in undergraduate and post graduate course development. He has been internal as well as external examiner of postgraduate thesis submissions. He was also local and technical chair of national and international conferences.

Ralf Klinkenberg

Ralf Klinkenberg holds Master of Science degrees in computer science with focus on machine learning, data mining, text mining, and predictive analytics from the Technical University of Dortmund in Germany and Missouri University of Science and Technology in the USA. He performed several years of research in these fields at both universities before initiating the RapidMiner open source data mining project in 2001, whose first version was called Yet Another Learning Environment (YALE). Ralf Klinkenberg founded this software project together with Dr. Ingo Mierswa and Dr. Simon Fischer. In 2006 he founded the company Rapid-I together with Ingo Mierswa. Rapid-I now is the company behind the open source software solution RapidMiner and its server version RapidAnalytics, providing these and further data analysis solutions, consulting, training, projects, implementations, support, and all kinds of related services. Ralf Klinkenberg has more than 15 years of experience in consulting and training large and small corporations and organizations in many different sectors how to best leverage data mining and RapidMiner based solutions for their needs. He performed data mining, text mining, web mining, and business analytics projects for companies like telecoms, banks, insurances, manufacturers, retailers, pharmaceutical com-

xxxvi

panies, healthcare, IT, aviation, automotive, and market research companies, utility and energy providers, as well as government organizations in many European and North American countries. He provided solutions for tasks like automated direct marketing campaign optimization, churn prediction and prevention, sales volume forecasting, automated online media monitoring and sentiment analysis to generate customer insights, market insights, and competitive intelligence, customer feedback analysis for product and service optimization, automated e-mail routing, fraud detection, preventive maintenance, machine failure prediction and prevention, manufacturing process optimization, quality and cost optimization, profit maximization, time series analysis and forecasting, critical event detection and prediction, and many other data mining and predictive analytics applications.

List of Contributors

Editors

- Markus Hofmann, Institute of Technology Blanchardstown, Ireland
- Ralf Klinkenberg, Rapid-I, Germany

Chapter Authors

- Ingo Mierswa, Rapid-I, Germany
- M. Fareed Akhtar, Fastonish, Australia
- Timm Euler, viadee IT-Consultancy, Münster/Köln (Cologne), Germany
- Matthew A. North, The College of Idaho, Caldwell, Idaho, USA
- Matej Mihelčić, Electrical Engineering, Mathematics and Computer Science, University of Twente, Netherlands; Rudjer Boskovic Institute, Zagreb, Croatia
- Matko Bošnjak, University of Porto, Porto, Portugal; Rudjer Boskovic Institute, Zagreb, Croatia
- Nino Antulov-Fantulin, Rudjer Boskovic Institute, Zagreb, Croatia
- Tomislav Šmuc, Rudjer Boskovic Institute, Zagreb, Croatia
- Milan Vukićević, Faculty of Organizational Sciences, University of Belgrade, Belgrade, Serbia
- Miloš Jovanović, Faculty of Organizational Sciences, University of Belgrade, Belgrade, Serbia
- Boris Delibašić, Faculty of Organizational Sciences, University of Belgrade, Belgrade, Serbia
- Milija Suknović, Faculty of Organizational Sciences, University of Belgrade, Belgrade, Serbia
- Andrew Chisholm, Institute of Technology, Blanchardstown, Dublin, Ireland
- Neil McGuigan, University of British Columbia, Sauder School of Business, Canada
- Eduarda Mendes Rodrigues, University of Porto, Porto, Portugal
- Luis Sarmento, Sapo.pt Portugal Telecom, Lisbon, Portugal
- Gurdal Ertek, Sabancı University, Istanbul, Turkey
- Dilek Tapucu, Sabancı University, Istanbul, Turkey

xxxviii

- Inanc Arin, Sabancı University, Istanbul, Turkey
- Tim Ruhe, TU Dortmund, Dortmund, Germany
- Katharina Morik, TU Dortmund, Dortmund, Germany
- Wolfgang Rhode, TU Dortmund, Dortmund, Germany
- Mertik Matej, Faculty of information study Novo mesto, Slovenia
- Palfy Miroslav, University Medical Centre Maribor, Slovenia
- Markus Muehlbacher, Department of Psychiatry and Psychotherapy, University Hospital of Erlangen-Nuremberg, Friedrich-Alexander-University Erlangen, Germany; Computer Chemistry Center, Friedrich-Alexander-University Erlangen, Germany
- Johannes Kornhuber, Department of Psychiatry and Psychotherapy, University Hospital of Erlangen-Nuremberg, Friedrich-Alexander-University Erlangen, Germany
- Radim Burget, Brno University of Technology, Czech Republic
- Václav Uher, Brno University of Technology, Czech Republic
- Jan Masek, Brno University of Technology, Czech Republic
- Marcin Blachnik, Silesian University of Technology, Department of Management and Informatics, Poland
- Miroslaw Kordos, University of Bielsko-Biala, Department of Mathematics and Computer Science, Bielsko-Biala, Poland
- Markus Goldstein, German Research Center for Artificial Intelligence, Kaiserslautern, Germany

Acknowledgments

A lot of people have contributed to make this book and the underlying open source software solutions RapidMiner and RapidAnalytics happen. We are thankful to all of you.

We would like to thank the contributing authors of this book, who shared their experience in the chapters and who thereby enable others to have a quick and successful data mining start with RapidMiner providing successful application examples and blueprints for the readers to tackle their data mining tasks and benefit from the strength of using RapidMiner and RapidAnalytics.

Many thanks to Dr. Brian Nolan, Head of Department of Informatics, Institute of Technology Blanchardstwon (ITB), for continuously supporting the relationship between the institute and Rapid-I.

The entire team of the Taylor & Francis Group were very professional, responsive and always helpful in guiding us through this project. Should any of you readers consider publishing a book then we can highly recommend this publisher.

Before there could be any thought about a book like this one, there needed to be the open source data mining software RapidMiner to whose success many contributed.

A special thanks goes to Prof. Dr. Katharina Morik, Head of the Artificial Intelligence Unit at the Technical University of Dortmund, Germany, for providing an introduction to and deep insights into machine learning (ML), data mining, text mining, artificial intelligence (AI), and natural language processing (NLP), for providing an environment that enabled the initiation of an open source data mining software project named YALE (Yet Another Learning Environment), which was later improved and renamed to RapidMiner. She supports the open source project RapidMiner, the company Rapid-I, which is behind the project, and its founders until today, long after they left the university. We appreciate the good cooperation and exchange of ideas with her research and teaching unit.

Another big thank you goes to Dr. Ingo Mierswa and Dr. Simon Fischer who started the open source project YALE and later RapidMiner together with Ralf Klinkenberg and who took turns in supervising the international development team of the project. Without their ideas, passion, commitment, and enormous development work, we would not have such a powerful and flexible open source data mining framework and solution today, available for everyone to use and gain from.

Dr. Ingo Mierswa and Ralf Klinkenberg are also the co-founders of Rapid-I, the data mining and business analytics company behind the open source software RapidMiner.

We are grateful to all who joined the RapidMiner and Rapid-I teams, especially to Sebastian Land, Helge Homburg, and Tobias Malbrecht who joined the RapidMiner team in its early days and contributed a lot to its development and who are strong members of the team until today. We are also thankful to all contributors from the early days, like Michael Wurst, Martin Scholz, and Timm Euler, as well as to the newer team members like Marius Helf, Nils-Christian Whler, Marcin Skirzynski, Venkatesh Umaashankar, Marco Bck, Dominik Halfkann, and to those who support the team in other roles, like Nadja Mierswa, Simone Horstmann, Christian Brabandt, Edin Klapic, Balzs Brny, Dietrich Niederlintner, Caroline Hahne, Miguel Bscher, Thilo Kamradt, Jannik Zappe, Kyle Goslin, and Assumpta Harvey. Open source projects grow strong with their community. We are thankful to all contributors to RapidMiner and RapidAnalyitcs and to all supporters of these open source projects. We are grateful not only for source code contributions, community support in the forum, bug reports and fixes, but also to those who spread the word with their blogs, videos, and words of mouth, especially to Thomas Ott,³ Neil McGuigan,⁴ Dr. Bala Deshpande,⁵ Prof. Dr. Bonnie Holub,⁶ Prof. Dr. Matthew North,⁷ Sheamus McGovern,⁸ Prof. Dr. David Wheismann, and many more.

Many bright minds have influenced our thoughts and inspired us with their ideas and valuable discussion. We would like to thank Prof. Dr. Thorsten Joachims, Prof. Dr. Hans-Ulrich Simon, Prof. Dr. Daniel St. Clair, Prof. Dr. Cihan Dagli, Prof. Dr. Tom Mitchell, and many others for widening our horizon and for many deep insights into the theoretical limits, enormous capabilities, and many practical applications of machine learning, data mining, data mining, text mining, statistical learning theory, and predictive analytics.

We would also like the many companies and organizations supporting the development of RapidMiner and RapidAnalytics by becoming customers of Rapid-I, including, for example, Sanofi, Daimler, Honda, Volkswagen, Miele, Siemens, Telekom Deutschland, T-Mobile International, mobilkom austria, Telenor, Nokia, Philips, Lufthansa, EADS, Salzgitter Mannesmann, ThyssenKrupp, Libri, Tchibo, KNV, PayPal, Intel, PepsiCo, GfK, Landesbank Berlin, E.ON, RWE, Axel Springer, 1&1, Schober, Schweizer Bundesbahn, FC Chelsea, The Cleveland Indians, and many more⁹.

With best regards and appreciation to all contributors,

Dr. Markus Hofmann, Institute of Technology Blanchardstown, Dublin, Ireland

Ralf Klinkenberg, Co-Founder of RapidMiner and Rapid-I, CBDO, Rapid-I, Dortmund, Germany

³http://www.NeuralMarketTrends.com/

⁴http://vancouverdata.blogspot.com/

⁵http://www.SimaFore.com/

⁶http://www.arclight.biz/uncategorized/rapidminer-introductory-tutorial-videos/

⁷http://docs.rapid-i.com/r/data-mining-for-the-masses

⁸http://www.cambridgecodeworks.com/wordpress/?author=2, Boston-Predictive-Analytics/

[?]author=2, http://www.meetup.com/

⁹http://rapid-i.com/content/view/8/56/
List of Figures

2.1	Welcome Perspective of RapidMiner.	20
2.2	Toolbar icons for perspectives	21
2.3	Design Perspective of RapidMiner.	22
2.4	A typical process in RapidMiner consists of several operators	23
2.5	Drag the Iris dataset into the process view in order to create a new operator	
	loading this dataset during process execution.	24
2.6	The probably most simple process which can be created with RapidMiner:	
	It just retrieves data from a repository and delivers it as a result to the user	~~
~ -	to allow for inspection.	25
2.7	Drag the operator named "Decision Tree" into your process	26
2.8	The complete process consisting of data loading and model creation	26
2.9	Press the Play icon in order to execute your process.	29
2.10	The decision tree described the rules for assigning the different classes to a	
	new plant	30
3.1	Workflow of the process.	39
3.2	Parameters of the Rename operator.	40
3.3	The parameters of the Numerical to Polynominal operator.	40
3.4	The parameters of the Numerical to Binominal operator.	41
3.5	The parameters of the Set Role Operator.	41
3.6	Subprocesses and Parameters of the Split Validation operator.	42
3.7	The change in accuracy of the model with the change in value of the	
0.1	parameter k.	43
4.1	Workflow of the process.	49
4.2	The dataset meta data information parameter of the Read CSV operator.	49
4.3	The parameters of the Split Data operator	51
5.1	Workflow of the process	57
5.2	Step 1 of Import Configuration Wizard of the Read CSV operator	58
5.3	Step 2 of Import Configuration Wizard of the Read CSV operator	59
5.0	Step 2 of Import Configuration Wizard of the Read CSV operator	59
5.5	Step 4 of Import Configuration Wizard of the Read CSV operator	60
5.6	Parameters of the Bename by Benlacing operator	60
5.0	The parameters of the Discretize By Binning operator	62
5.8	Training and Testing subprocesses of the X-Validation operator	62
5.9	Effect of discretization and filtering on the accuracy of the Naïve Bayes	02
0.0	model.	63
6.1	Workflow of the process.	68
6.2	Step 1 of import Configuration Wizard of the Read Excel operator	68
6.3	Step 2 of import Configuration Wizard of the Read Excel operator	69

xlii

$ \begin{array}{r} 6.4 \\ 6.5 \\ 6.6 \\ 6.7 \\ 6.8 \\ \end{array} $	Step 3 of Import Configuration Wizard of the Read Excel operator.Step 4 of import configuration Wizard of the Read Excel operator.The parameters of the Split Data operator.First row of labeled ExampleSet.Distribution table.	70 70 71 72 73
7.1 7.2 7.3 7.4	The source data schema	80 82 83
7.5	account data	83
	ing account data.	84
$7.6 \\ 7.7$	Indicator attributes at the account level	85
7 0	erator)	85
1.0 7.0	Ine substream in Chr_ol_Createwining rable for processing transaction data.	00
7.9	A BapidMiner contingency table	80
7.10	A RapidMiner process for creating lift charts	90
7.12	A lift chart created with BapidMiner	91
7.12	A ROC chart created with RapidMiner on real data.	92
8.1	A simplified relational model of supermarket's database	98
8.2	An example customer record from a supermarket's database	99
8.3	An example product record from a supermarket's database	99
8.4	Example of product categories in a supermarket's database	99
8.5	Examples of receipt data in a supermarket's database	100
8.6	Products connected to their corresponding receipts in a supermarket's database.	100
8.7	Query results indicating which product categories are found on each receipt.	101
8.8	Query results from an expanded dataset, ready for further analysis in Rapid-	101
89	Supermarket data, extracted from the relational database in CSV format	101
8 10	Connecting to a CSV file in BapidMiner	102
8.11	Connecting to an unmapped network location.	103
8.12	Setting the column delimiters.	103
8.13	Setting the first row in our CSV import as our attribute names	104
8.14	Assignment of data types to each attribute.	104
8.15	Saving the imported dataset.	105
8.16	Using the Read CSV operator.	106
8.17	An example of binominal receipt data in a MySQL database table	106
8.18	Using the Read Database operator.	107
8.19	Configuring Read Database parameters.	107
8.20	Setting up the process window to build a model	108
8.21	Results Perspective after running a model with only a Retrieve object in it.	109
8.22	Addition of a Select Attributes operator to the model	110
8.23	Configuring the Select Attributes operator	110
8.24	A Sample operator configured to take a 50% sample of the data	111
8.25	Using the Filter Examples operator to remove any records containing miss-	
	ing values in any of the attributes in the dataset	111

8.26	Inconsistent data that will complicate an Association Rule model	112
8.27	Inclusion of a Declare Missing Value operator, removing any '2' values	112
8.28	In Results View, the '2' values have been changed to 'missing' (?)	112
8.29	Addition of a Replace Missing Values operator, changing missing values to	
	'1'	113
8.30	The addition of FP-Growth to our data mining process.	113
8.31	Item sets generated by FP-Growth.	114
8.32	Create Association Rules in the main process of our data mining model.	114
8.33	Results of the Associate Rules model at 80% confidence.	115
8.34	The Association Rule model with a 75% confidence percentage	115
8.35	Graph View of our Association Rules in Results Perspective	116
9.1	An example of an item recommendation workflow	130
9.2	Measuring performance of a recommendation model	131
9.3	Item recommendation online update workflow	132
9.4	Item attribute recommendation workflow.	133
9.5	Similarity-based content recommendation workflow. Similarity is based on	
	textual similarity between consumed and new items	134
9.6	Text analysis workflow for VideoLectures Case study	136
9.7	Hybrid recommender system for VideoLectures Case study	136
9.8	Simple recommendation web service workflow (upper image) and the content	
	of the ProcessInput operator (lower image).	138
9.9	Settings and test of the simple recommender web service	139
9.10	Architecture of recommendation web engine	140
10.1	Process for automatic evaluation of classification algorithms	148
10.2	Discretization of the output attribute	149
10.3	Assigning of "label" Role to the "Average_Grade" attribute.	149
10.4	Separation of the data by study programs	149
10.5	Separation of the data by study programs	150
10.6	"Simple Validation" and logging operators	150
10.7	"Simple Validation" inner operators.	151
10.8	Classification algorithms as inner operators of the "Select Sub-process"	151
10.9	Classification algorithms as inner operators of the "Select Sub-process"	151
10.10	"Wrapper Split Validation" operator	152
10.11	Inner operators of "Wrapper Split Validation".	153
10.12	"X-Validation" for the estimation of algorithm performance with feature	
	selection.	153
		101
11.1	3-D scatter plot of artificial data showing 8 clusters	164
11.2	$Raw E-coli data. \ldots \ldots$	165
11.3	Overview of the RapidMiner process.	167
11.4	Detail of Generate Clusters	170
11.5	Generate ground truth measures detail	171
11.6	First few rows of $k = 2$ k-means clustering merged with original clusters for	
	artificial data.	172
11.7	Operator chain to calculate external validity measures using R	172
11.8	Example set returned by R process	173
11.9	Internal validity measure calculation	174
11.10	Process connection for artificial data	175
11.11	Process connection for <i>E-coli</i> data	176

xliii

 \mathbf{x} liv

$11.12 \\ 11.13 \\ 11.14$	Graph of Rand, Jaccard, Fowlkes-Mallow, and adjusted Rand indexes Internal validity measures as a function of k for artificial data Adjusted Rand index validity measure between different clusterings for ar-	177 178
11.15	tificial data	179
$11.16 \\ 11.17$	<i>E-coli</i> data	180 180 181
$12.1 \\ 12.2 \\ 12.3 \\ 12.4 \\ 12.5 \\ 12.6 \\ 12.7 \\ 12.8 \\ 12.9 \\ 12.10$	Main process for automatic evaluation of clustering algorithms Attribute selection menu	188 189 189 190 190 191 191 192 193
$13.1 \\ 13.2 \\ 13.3 \\ 13.4 \\ 13.5 \\ 13.6 \\ 13.7$	Data Import Wizard.The basic Text Analysis Process.A table of word frequencies.Balancing the Sample by the Label.The Text Classification Process.Cross-Validation.Inside the X-Validation operator.	202 204 205 206 208 209 209
$\begin{array}{c} 14.1 \\ 14.2 \end{array}$	Examples of different alphabets	216
$14.3 \\ 14.4$	ences in the use of distinct letters per language	217 224
14.5	words in languages	227
14.6	(lower image)	228
14.7	The workflow for the language identification system based on character n- grams is depicted in the upper image	229 231
14.8	Estimated performance of the n-grams method with 10-fold cross-validation (upper image) and performance of the method on the test set (lower image).	231
14.9	List of the top 10 most frequent n-grams in the language identification dataset, obtained with the character n-gram workflow.	232
14.10	Workflow for the language identification system based on n-grams profile similarity	233
14.11	Performance of the n-gram similarity method on the test set	233
14.12	The workflow of the web page classification application	235
14.13	Exporting the web page classification workflow as a service.	236
14.14	Using the web page classification service	237

15.1	Operators for Process01 and the Parameters for Process Documents from Files	244
15.2	Operators within the Process Documents from Files nested operator	$\frac{245}{245}$
15.2	Parameters for the operators within the Process Documents from Files oper-	- 10
10.0	ator	246
$15\ 4$	Parameters for the operators in Process01	$\frac{2}{246}$
15.1	Configuring LocalRepository	$\frac{2}{247}$
15.6	Operators in Process02	248
15.0	Operators within the Process Documents from Files nested operator and the	- 10
10.1	Parameters for the Generate n-Grams (Terms) operator.	249
15.8	Parameters for the operators in Process02.	$\frac{1}{249}$
15.9	Opening and running Process01.	251
15.10	Dialog box alerting the switch to the result perspective	251
15 11	Result Overview for Process01 results	251
15.12	Specifying the data source text directories for Process01	$\frac{101}{252}$
15 13	Result Overview for Process01 results	253
15.10	Wordlist generated by Process01	$\frac{100}{253}$
15 15	Meta Data View for the ExampleSet generated by Process01	$\frac{100}{254}$
15.16	Data View for the ExampleSet generated by Process01	$\frac{101}{254}$
15.10	Table View for the AssociationRules generated by Process01	255
15.17	Filtering rules in the Graph View for the AssociationRules	$\frac{200}{256}$
15.10	Document Occurrences of the words in the Word list	256 256
15 20	Result Overview for Process02 results	$\frac{100}{258}$
15.20	Meta Data View for the ExampleSet generated by Process02 including the	200
10.21	n-Grams	259
15 22	Data View for the ExampleSet generated by Process02 and the relative con-	200
10.22	currence frequency of the word absolut in hotel 73943.txt	259
15 23	Centroid Table for the Cluster Model generated by Process02 displaying the	_00
10.20	average frequency of each word in each cluster	260
		200
16.1	Overview over the wide field of astroparticle physics	266
16.2	Schematic view of the detection principle. A cone of Cherenkov light is	
	emitted along the track of the particle. Track parameters, e.g., the angle of	
	the incident neutrino, can be reconstructed from the light pattern recorded	
	by the optical modules	268
16.3	View of the Feature Selection Extension after a successful installation.	269
16.4	View of the Feature Selection Setup.	270
16.5	Inner process of the LOOP PARAMETERS operator.	271
16.6	Inner process of the WRAPPER X-VALIDATION	272
16.7	Editing the settings of the LOOP PARAMETERS operator.	274
16.8	Looping over a list of values instead of a grid.	275
16.9	Setup for the evaluation of the Feature Selection Stability.	276
16.10	Settings of the LOOPPARAMETERS operator used for the stability evaluation.	277
16.11	Inner process, FSSV, of the LOOP PARAMETERS operator.	278
16.12	Settings of the LOG operator.	278
16.13	Weighting operator placed inside the FEATURE SELECTION STABILITY VALIDA-	
	TION	279
16.14	Setup of the complete learning process using a 5-fold cross validation.	279
16.15	Settings for the GENERATE ATTRIBUTES operator.	280
16.16	Settings for the REMAP BINOMINALS operator.	280
16.17	Setup of the training process inside the CROSS VALIDATION operator.	281

	•
x	V1

$\begin{array}{c} 16.18\\ 16.19 \end{array}$	Settings for the W-RANDOM FOREST operator	$\begin{array}{c} 282 \\ 283 \end{array}$
16.20	Overview over both operators included in the subprocess	284
16.21	Settings of the EXTRACT MACRO operator	284
16.22	Setting to create a new attribute using a macro	285
1 - 1		205
17.1	Meta Data View on CTS dataset.	295
17.2	Simple data visualisation: Deviation Plot.	296
17.3	Simple data visualisation: Multiple Series Plot	297
17.4	Operators used in pre-processing	298
17.5	Select attributes dialog	299
17.6	Whole KDD process.	301
17.7	Learning in X-Validation operator.	301
17.8	Parameters for ANN	302
17.9	Preprocessed example set.	303
17.10	Neural Network Model.	304
17.11	Performance vector	304
17.12	Statistical process for hypothesis testing	305
17.13	Reconstruction subprocess for ANOVA	306
17.14	Renaming attributes in pre-processing	307
17.15	ANOVA matrix.	308
17.16	ExampleSet after Aggregate operator	308
17.17	Plot for temperature differences	309
17.18	Meta Data View of the diabetes dataset.	310
17.19	Descriptive data mining—searching for clusters	311
17.20	Diabetes—attribute selection.	312
17.21	ResultViews for diabetes	313
17.22	Various trees.	317
17.23	Hypothesis testing process for diabetes.	318
17.24	Result of the hypothesis test.	318
	v I	
18.1	After successful installation the operator tree shows the <i>Chemistry</i> folder, which includes all 10 operators of the PaDEL extension, after installation.	322
18.2	The workflow for the complete import, preprocessing, and descriptor cal-	
	culation process. The operators 7, 8, and 9 are provided by the PaDEL	
	extension.	327
18.3	Example workflow to generate a simple quantitative model for the prediction	
	of aqueous solubility.	327
18.4	This scatterplot shows the experimental solubility (S) and the prediction of	
	S	328
19.1	This workflow outlines the import and preprocessing of the example dataset	334
19.2	The figure illustrates the complete feature selection process	336
19.3	The workflow illustrates the nested process of the "Optimize Selection"	
10.4	operator performing a beam search.	337
19.4	A validation consists of a nested training and testing process	338
19.5	The Y-randomization workflow generates models based on randomly shuf-	
	tled label values. These models should not be as predictive as the original	
	model.	339
19.6	The plotted evaluation of the beam search shows that there is an increase	_
	of the accuracy up to 4 attributes	340

xlvii

20.1	How to enable expert mode in RapidMiner	348
20.2	Scheme of image mining	349
20.3	Example: Read image, blur and save result.	350
20.4	Example: Iterate over set of images	351
20.5	Example: Color component extraction.	353
20.6	Example: Combining images.	354
20.7	Example: How to remove objects according to their shape	355
20.8	Principle of feature extraction.	356
20.9	Example: Trainable segmentation.	356
20.10	Example: Trainable segmentation.	357
20.11	Example: How to remove objects according to their shape	358
20.12	Example: How to remove objects according to their shape	359
20.13	Example: Image classification.	360
	I GOLDAN AND	
21.1	Example of two image classes	364
21.2	Process tree	365
21.3	Choosing the class	365
21.4	Image loading settings	366
21.5	OBCF operator function.	366
21.6	Pre-processed image by Border / Interior Classification operator.	367
21.7	Operators for global features extraction.	368
21.8	Haar-like features	369
21.9	Ultrasound artery image	369
21.10	Samples for artery detection.	370
21.11	Read Image Set operator settings	370
21.12	Viola-Jones operator settings	371
21.13	Process for pattern detection	371
21 14	Process overview	372
21 15	Example of segmentation	373
		0.0
22.1	Equivalent examples of applications of the Instance Selection operators (a)	
	use of internal k-NN model (b) use of selected instances to train external	
	k-NN model.	380
22.2	Validation of the accuracy and compression of the <i>GE Selection</i> operator	384
22.3	Configuration of <i>Log</i> operator, and the table view of the results	386
22.4	Selecting the most useful instance selection operator. Validation inner pro-	
	cess configuration.	387
22.5	Selecting the most useful instance selection operator Main process	387
22.6	Selecting the most useful instance selection operator Validation inner pro-	
	cess configuration.	388
22.7	Processing many datasets at once.	390
22.8	Test of the outlier elimination task.	391
22.9	Process configuration used to compare different instance selection scenarios	
	(1)	393
22.10	Process configuration used to compare different instance selection scenarios	
0	(2)	394
22.11	Comparison of three different combinations of instance selection algorithms	
	(CNN, ENN+CN, ENN+CNN+RNG, accuracy as a function of Compres-	
	sion plot, and configuration settings required to generate it.	396
22.12	Process of prototype construction based on clustering and re-labeling	397
		501

xlviii

22.13	Process of prototype construction based on clustering of each class and combing obtained cluster centres.	397
22.14	Process of prototype construction based on clustering of each class and combing obtained cluster control	208
22.15	Process tree of prototype construction based on the LVO algorithm	300
22.10 22.16	Accelerating the comparison of three LVO algorithms using subresults	000
22.10	caching the comparison of three LVQ algorithms using subresults	400
22 17	Evaluation process of the time complexity of the instance selection methods	100
22.11	(1)	402
22.18	Evaluation process of the time complexity of the instance selection methods	102
		403
22.19	Accuracy, training time, and testing time in a function of the sample size	
	using different instance selection methods followed by the SVM (no instance	
	selection, CNN, ENN, and ENN+CNN).	405
22.20	Select Subprocess operators configuration used to evaluate the time com-	
	plexity of the instance construction methods.	406
22.21	Accuracy, training time, and testing time in a function the sample size using	
	SVM and two instance construction algorithms LVQ and K-means	407
23.1	The result of data cleansing using the <i>Detect Outlier (Distances)</i> operator	
	on the <i>iris</i> dataset	411
23.2	A point anomaly detection problem	414
23.3	A contextual anomaly detection problem: The average monthly temperature	414
09.4	In Germany from Jan 2001 to Dec 2010	414
23.4	A contextual anomaly detection problem having two contexts: 1 ne time and the user ID	415
02 E	The user ID	410
20.0 22.6	The process for generating an artificial 2D anomaly detection data set	410
23.0	randomly generated as normal data using a mixture of Caussians	/18
23.7	The parameter settings of the Man operator A regular expression is used	410
20.1	to map the cluster names in the label attribute to a single label entitled	
	"normal".	418
23.8	The scatter plot of the generated artificial 2D dataset for a unsupervised	110
	anomaly detection process. Blue color indicates normal data instances, red	
	color the sampled outliers.	419
23.9	The result of the global k-NN anomaly detection.	420
23.10	A zoomed view of the c_2 cluster of Figure 23.2 illustrating the local outlier	
	p_4 detected by LOF. The circles represent the LRDs for $k = 3$. The LOF	
	score is the ratio of the light grey circle to the average of the dashed circles.	421
23.11	Results of a straight line dataset with two outliers using LOF. The top two	
	outlier scores are marked with red color	422
23.12	Results of a straight line dataset with two outliers using COF. The top two	
	outlier scores are marked with red color	422
23.13	A sample process for applying a cluster-based anomaly detection algorithm	
	on the data. On the right side, the parameter settings for CBLOF on the	
00.1.1	artificial dataset are shown.	423
23.14	The result of the LDCOF anomaly detection algorithm using X-means clus-	
	tering. The bubble size indicates the outlier score and the color indicates	49.4
92 1E	The proprocessing of the NRA data: Personal attributes concerting on	424
29.19	appropriate ID, and removing missing values	196
	appropriate ID, and removing informs values	440

3.2	1 7 7 7
× 1	II V
	пл

23.16 23.17 23.18	Processing the NBA data	427 429
	dataset	430
23.19	A complex process for semi-supervised cluster-based anomaly detection	431
23.20	The inner loop of the process	433
23.21	Using the <i>Execute Script</i> operator for computing anomaly scores based on	
	the distances to the nearest cluster center $\ldots \ldots \ldots \ldots \ldots \ldots \ldots$	433
24.1	The first setting for the evaluation of learning algorithms	443
24.2	Looping through different datasets	444
24.3	Parameter setting for the Loop Files operator	444
24.4	Selecting one of the learning algorithms.	446
24.5	Parameter for the <i>Select Sub-process</i> operator	447
24.6	Loops for datasets and algorithms combined	447
24.7	Operators for logging the results	448
24.8	Logging results.	449
24.9	Post-processing of the results.	449
24.10	Results in the final format	450
24.11	Statistical analysis of the results.	451
24.12	Results of the statistical t-test for Australian (upper) and German (lower).	452
24.13	Extraction of dataset properties for meta-learning.	453
24.14	Log output with additional dataset properties	454

List of Tables

1.1	A set of examples with values for all attributes together with label values.	16
$6.1 \\ 6.2$	Posterior probabilities	72 73
9.1 9.2 9.3	Recommendation operators supported by the Recommender Extension Summary of the requirements and effects of the recommendation operators. An example of an AML and a related DAT file for item recommendation	122 123
5.0	operators.	124
9.4 9.5	Implicit feedbacks in movie recommender system. "1" denotes user con- sumed an item, "0" denotes he/she did not	128 130
$10.1 \\ 10.2 \\ 10.3$	Description of I/O variables used in this study	148 154 154
$\begin{array}{c} 11.1 \\ 11.2 \end{array}$	Process sections	168 169
12.1	Description of attributes for clustering students	187
13.1	Fictional historical data on car theft.	207
14.1	Lists of the 20 most frequent words in Dutch, Romanian, and Swedish	226
$17.1 \\ 17.2 \\ 17.3 \\ 17.4$	Predictive functions and associated algorithms for medical data mining Descriptive functions and associated algorithms for medical data mining . Selection of DM algorithm considering type of (medical) data	293 294 300 300
18.1	Operator descriptions.	323
19.1	Summary of the results achieved using the workflow from this use case	340
22.1 22.2 22.3	List of datasets used in the experiments	378 385 392
23.1	Top-10 anomalies found in the NBA dataset using the k-NN global outlier detection algorithm.	428

23.2	The confusion matrix of the ground truth and the semi-supervised predic-	
	tions for the shuttle dataset using a one-class SVM	430
23.3	The confusion matrix of the ground truth and the semi-supervised predic-	
	tions for the shuttle dataset based on our simple cluster-based anomaly	
	score	431

lii

Part I

Introduction to Data Mining and RapidMiner

Chapter 1

What This Book is About and What It is Not

Ingo Mierswa

Rapid-I, Dortmund, Germany

1.1	Introduction				
1.2	Coinci	Coincidence or Not?			
1.3	Applie	cations of Data Mining	7		
	1.3.1	Financial Services	7		
	1.3.2	Retail and Consumer Products	8		
	1.3.3	Telecommunications and Media	9		
	1.3.4	Manufacturing, Construction, and Electronics	10		
1.4	Funda	mental Terms	11		
	1.4.1	Attributes and Target Attributes	11		
	1.4.2	Concepts and Examples	13		
	1.4.3	Attribute Roles	14		
	1.4.4	Value Types	14		
	1.4.5	Data and Meta Data	15		
	1.4.6	Modeling	16		

1.1 Introduction

Today, analytics is a very important topic and affects practically all levels in modern organizations. Analytics is also used by many data-driven researchers. Data is collected and analyzed, and the results of this analytical work either prove our hypothesis or delivers new insights.

When we talk about analytics in this book we are referring to what many call "advanced analytics". This field includes technologies known from statistics as well as from computer science. Isn't this just statistics done by computers then? By far not! Statistics often deal with the question if a hypothesis can be proven by a statistical test using a small but representative data sample. Although this is of biggest importance, it is even more useful to mix these ideas with algorithms from computer science to make sure that the methods we are talking about will be able to scale up and analyze even the largest datasets.

And I see another distinction: traditional statistics often requires that the analyst is creating a model or hypothesis right at the start of the analysis process. After creating such a model, parameters of the model are estimated or the applicability of this model is proven by means of a statistical test. Maybe this is because I am lazy, but I actually don't like this idea too much: Why should I manually do the work a computer is perfectly able to deal with itself? In this sense of manual analysis, statistical analysis is much more connected to online analytical processing (OLAP) than with "advanced analytics": In the OLAP world, people try to drill through their data to find the interesting patterns and the reasons in deeper data levels themselves. Fine, but again I think this is the wrong approach for solving

the underlying problem for mainly two reasons: Firstly, people tend to see only what they are looking for. Most analysts have some expectations before they start and try to work as long as necessary on the data to prove their point. Secondly, OLAP is again a pretty tedious work I personally believe a computer is much better suited for. Did I already mention that I am pretty lazy in this respect? I prefer actually to describe myself as "unbiased" and "efficient".

Don't get me wrong: statistics and OLAP offer very important methods necessary for many day-to-day business cases and I myself am half computer scientist and half statistician. However, if you mix the methods described above with algorithms from computer science to scale those methods up to larger datasets and also throw in some ideas from artificial intelligence, especially from the field of machine learning, I personally think that much more interesting possibilities can arise. This is actually nothing new and has been an important field for research during the last decades. The methods and algorithms which have been developed over this time have actually formed complete new research fields under the names data mining, predictive analytics, or pattern detection. And one of the most amazing developments is a collection of methods which can not only be used on structured data, i.e., on tables, but also on unstructured data like texts or images. This has been the underlying motivation for fields like text mining, image mining, or audio mining.

Most recently a new buzzword has been used a lot: Big Data. Well, most recently means in the years 2012 and following, so if you are reading this book in 2092 you might want to use this as a historical lecture. OK, back to big data: What is so special about that? If you ask me, not that much from the point of view of a data miner. Big data is an umbrella term for many ideas and technologies, but the underlying point of all those things is that big data should be about infrastructure and methods for collecting, retrieving, and analyzing very large datasets which might be of structured, unstructured, or polystructured nature. Well, if you have read the paragraphs above you will certainly agree that this sounds actually like a perfect description of "data mining". As of 2013, the big data market is still in its early days and most people are worrying about data infrastructures. But this will change as soon as organizations do understand that the mere collection of data is worth nothing. It is the analysis of data which delivers new insights, explains underlying patterns, or creates models which can be extrapolated in order to predict the most likely future.

Reading this book might therefore be a very good idea for learning more about where and how data mining can be used to deliver those insights. It might also just serve you for your personal career—remember, the big data market will slowly but surely move into the direction of analytics in the future. Whatever the reason is, I hope that you will learn more about the use cases discussed in this book so that you are able to transfer them to your own business problems. RapidMiner is an excellent and very flexible tool for re-using those use cases and adapt them to your concrete problems. Have fun using it!

In the following, I would like to discuss what data mining actually is and to which types of problems it can be applied. At the same time I will give you an introduction to the most important terms. Whether you are already an experienced data mining expert or not, this chapter is worth reading in order for you to know and have a command of the terms used both here in this book as well as in RapidMiner.

1.2 Coincidence or Not?

Before we get properly started, let us try a small experiment:

- 1. Think of a number between 1 and 10.
- 2. Multiply this number by 9.
- 3. Work out the checksum of the result, i.e., the sum of the numbers.
- 4. Multiply the result by 4.
- 5. Divide the result by 3.
- 6. Deduct 10.

The result is 2.

Do you believe in coincidence? As an analyst you will probably learn to answer this question in the negative or even do so already. Let us take for example what is probably the simplest random event you could imagine, i.e., the toss of a coin. "Ah", you may think, "but that is a random event and nobody can predict which side of the coin will be showing after it is tossed". That may be correct, but the fact that nobody can predict it does in no way mean that it is impossible in principle. If all influence factors such as the throwing speed and rotation angle, material properties of the coin and those of the ground, mass distributions and even the strength and direction of the wind were all known exactly, then we would be quite able, with some time and effort, to predict the result of such a coin toss. The physical formulas for this are all known in any case.

We shall now look at another scenario, only this time we can predict the outcome of the situation: A glass will break if it falls from a certain height onto a certain type of ground. We even know in the fractions of a second when the glass is falling that there will be broken glass. How are we able to achieve this rather amazing feat? We have never seen the glass which is falling in this instant break before and the physical formulas that describe the breakage of glass are a complete mystery for most of us at least. Of course, the glass may stay intact "by chance" in individual cases, but this is not likely. For what it's worth, the glass not breaking would be just as non-coincidental, since this result also follows physical laws. For example, the energy of the impact is transferred to the ground better in this case. So how do we humans know what exactly will happen next in some cases and in other cases, for example that of the toss of a coin, what will not?

The most frequent explanation used by laymen in this case is the description of the one scenario as "coincidental" and the other as "non-coincidental". We shall not go into the interesting yet nonetheless rather philosophical discussions on this topic, but we are putting forward the following thesis:

The vast majority of processes in our perceptible environment are not a result of coincidences. The reason for our inability to describe and extrapolate the processes precisely is rather down to the fact that we are not able to recognize or measure the necessary influence factors or correlate them.

In the case of the falling glass, we quickly recognized the most important characteristics, such as the material, falling height, and nature of the ground, and can already estimate, in the shortest time, the probability of the glass breaking by analogy reasoning from similar experiences. However, it is just that we cannot do with the toss of a coin. We can watch as many tosses of a coin as we like; we will never manage to recognize the necessary factors fast enough and extrapolate them accordingly in the case of a random throw.

So what were we doing in our heads when we made the prediction for the state of the glass after the impact? We measured the characteristics of this event. You could also say that we collected data describing the fall of the glass. We then reasoned very quickly by analogy, i.e., we made a comparison with earlier falling glasses, cups, porcelain figurines,

6

or similar articles based on a similarity measure. Two things are necessary for this: firstly, we need to also have the data of earlier events available and secondly, we need to be aware of how a similarity between the current and past data is defined at all. Ultimately we are able to make an estimation or prediction by having looked at the most similar events that have already taken place, for example. Did the falling article break in these cases or not? We must first find the events with the greatest similarity, which represents a kind of optimization. We use the term "optimization" here, since it is actually unimportant whether we are now maximizing a similarity or the sales figures of one enterprise or any other—the measurement concerned, in this case similarity, is always optimized. The analogy reasoning described then tells us that the majority of glasses we have already looked at broke and this very estimation then becomes our prediction. This may sound complicated, but this kind of analogy reasoning is basically the foundation for almost every human learning process and is done at a staggering speed.

The interesting thing about this is that we have just been acting as a human data mining method, since data analysis usually involves matters such as the representation of events or conditions and the data resulting from this, the definition of events' similarities and of the optimization of these similarities.

However, the described procedure of analogy reasoning is not possible with the toss of a coin: It is usually insufficient at the first step and the data for factors such as material properties or ground unevenness cannot be recorded. Therefore, we cannot have these ready for later analogy reasoning. This does in no way mean, however, that the event of a coin toss is coincidental, but merely shows that we humans are not able to measure these influence factors and describe the process. In other cases we may be quite able to measure the influence factors, but we are not able to correlate these purposefully, meaning that computing similarity or even describing the processes is impossible for us.

It is by no means the case that analogy reasoning is the only way of deducing forecasts for new situations from already known information. If the observer of a falling glass is asked how he knows that the glass will break, then the answer will often include things like "every time I have seen a glass fall from a height of more than 1.5 meters it has broken". There are two interesting points here: The relation to past experiences using the term "always" as well as the deduction of a rule from these experiences:

If the falling article is made of glass and the falling height is more than 1.5 meters, then the glass will break.

The introduction of a threshold value like 1.5 meters is a fascinating aspect of this rule formation. For although not every glass will break immediately if greater heights are used and will not necessarily remain intact in the case of lower heights, introducing this threshold value transforms the rule into a rule of thumb, which may not always, but will mostly lead to a correct estimate of the situation. Instead of therefore reasoning by analogy straight away, one could now use this rule of thumb and would soon reach a decision as to the most probable future of the falling article. Analogy reasoning and the creation of rules are two first examples of how humans, and also data mining methods, are able to anticipate the outcome of new and unknown situations.

Our description of what goes on in our brains and also in most data mining methods on the computer reveals yet another interesting insight: The analogy reasoning described does at no time require the knowledge of any physical formula to say why the glass will now break. The same applies for the rule of thumb described above. So even without knowing the complete (physical) description of a process, we and the data mining method are equally able to generate an estimation of situations or even predictions. Not only was the causal relationship itself not described here, but even the data acquisition was merely superficial and rough and only a few factors such as the material of the falling article (glass) and the falling height (approx. 2m) were indicated, and relatively inaccurately at that. Causal chains therefore exist whether we know them or not. In the latter case, we are often inclined to refer to them as coincidental. And it is equally amazing that describing the further course is possible even for an unknown causal chain, and even in situations where the past facts are incomplete and only described inaccurately.

This section has given you an idea of the kind of problems we wish to address in this book. We will be dealing with numerous influence factors, some of which can only be measured insufficiently or not at all. At the same time, there are often so many of these factors that we risk losing track. In addition, we also have to deal with the events which have already taken place, which we wish to use for modeling and the number of which easily goes into millions or billions. Last but not least, we must ask ourselves whether describing the process is the goal or whether analogy reasoning is already sufficient to make a prediction. And in addition, this must all take place in a dynamic environment under constantly changing conditions—and preferably as soon as possible. Impossible for humans? Correct. But not impossible for data mining methods.

1.3 Applications of Data Mining

indexData Mining Applications Before we start to discuss the fundamental terms of data mining and predictive analytics, I would like to give you a feeling about the possible application fields for those techniques. There are literally hundreds of possible applications for data mining and predictive analytics across basically all verticals and horizontals you can think of. The following selection is therefore by no means complete; it should only give you some ideas where these technologies already have been successfully applied. The applications described in the following are grouped along selected verticals where they are most often used, but make sure to check them all since there is no reason why you should not use "fraud detection" also in retail instead of financial services.

1.3.1 Financial Services

- *Fraud detection:* Fraud detection is often used in the financial service industry but not only there. The basic idea here is to detect fraudulent activities among a very large set of transactions with methods from predictive analytics. Possible approaches include the detection of outliers or the modeling of "normal" cases against "fraudulent" cases and using this model for new transactions to check if they fall into the fraud segment.
- *Churn prevention:* Assume your insurance company has a contract with one of your customers. The optimal case for your insurance would that this contract stays and you keep an ongoing relationship with your customer. Sadly, some customers decide to quit the contract for a wide variety of reasons and you would like to know in advance for which customers this will happen with the highest likelihood in the near future. This is exactly the idea behind churn prevention: create predictive models calculating the probability that customers are likely to quit contracts soon so you can be pro-active, engage with them, offer incentives, etc. Besides the financial industry, churn prevention can also be often found in the retail industry, e-commerce, or telecommunications among others.
- Sentiment analysis: Sentiment analysis is not at all connected to the financial industry but we have seen it very often here lately. You can also find sentiment analysis

in industries like consumer goods, retail, telecommunications, and life sciences a lot. The idea behind sentiment analysis is to connect to thousands of online sources in the web, collect statements about your brands or products, and analyze them by means of text analytics with respect to their tonality or sentiment. You can identify how sentiment changes over time and measure the success of marketing or PR by this or you can get new insights about how to improve your products. If you combine this with network analysis you can even detect key opinion leaders and see how they influence their peer groups.

- **Trading analytics:** If you are trading, building portfolios, or preparing deals, the natural idea would be to calculate the success rate of your decisions with help of predictive analytics. In general, you could analyze potential trade opportunities by looking at market data or inspect markets for trading activity which showing an emerging trend. Lately, many analysts combine more traditional methods like time series analysis with behavioral trading algorithms or even sentiment analysis.
- *Risk management:* This is another example done in the financial services industry which can also be applied to many other industries as well, especially when it comes to supply chain management, for example in manufacturing, or to logistics and transportation. Data mining and predictive analytics can be used for solving multiple problems connected to risk management, including error detection and quantification, unit reviews or internal audits, detecting fraud (see above), identifying suppliers with highest probability of failure, quantifying payment risks, and credit scoring.

1.3.2 Retail and Consumer Products

- Customer segmentation, channel selection, and next best action: In retail you can easily find most applications of data mining for sales and marketing. A typical approach is to use data about customers, including descriptions of the customers' purchase and transaction history, for creating segments of customers, for example based on classification but more often based on clustering techniques. Those clusters can build data-driven segments much more optimized than the gut-driven A-B-C segments which even today can often be seen in practice. The assignment of customers to segments is an important prerequisite for further analysis, for example, for selecting customers for specific sales or marketing channels or for predicting which is the optimal next best action to approach those customers or leads.
- Direct marketing: The bread-and-butter business for data mining and one of the stories how it all has begun. The idea behind direct marketing is to assign costs for different types of actions: If I contact a lead and he does not react, this costs A (for the contact efforts, etc.). If he does, I will get some gain B (gains are essentially negative costs). If I don't contact the lead and he would not have reacted anyway, this saves me cost A above. But if I decide to not contact the lead and he or she would have purchases, this might cause a huge loss C. So the whole point is to identify those leads with the highest probabilities for conversion for a certain marketing campaign so that you only would contact the most likely cases up to a limit where the contacting costs would no longer lead to revenue high enough to compensate the costs. Even with the advent of e-mail, direct marketing is necessary, the "cost" here might be that recipients are tired of getting too much spam and might opt out instead so we are using the lead.
- Recommendations, cross-selling and up-selling: Another one of the big success

stories of data mining. Everybody who has already purchased a book at Amazon already came across the results of so-called recommender systems: "People who bought THIS book also purchased THAT one." At first sight, this problem type might not look too complicated: for each item, just search for those which have been purchased frequently together with the first one. The problem comes with the high number of available items and the high number of transactions typically available in retail and e-commerce. It is literally impossible to make those calculations for all combinations which might happen and so we need algorithms guiding our search to the most promising directions. We call those promising combinations "frequent item sets" and after we found those sets, we might want to recommend other items from those sets if the first one is added to the cart. This approach is called *cross-selling* and might complement or even replace traditional cross-selling approaches based on manual rules. Loosely connected to this is *up-selling* where we try to identify customers who are likely to purchase a higher-valued product or a larger quantity.

• **Customer lifetime value**: Traditional systems for business intelligence based on OLAP approaches are great for answering questions like "who are the customers who bought most so far" or "what are the revenues created with my top 10 customers in the past". Although this is with no doubt important information, it unfortunately is also only reflecting the past. Previously good customers might change to another supplier or drop out for other reasons, and the fact that a customer has created much revenue so far does not mean this will continue also in the future. Instead of analyzing the historical customer value, many organizations now turn to predict how customers will develop in the future and what their total customer lifetime value will be to guide their sales and marketing efforts. Predictive analytics methods help to identify typical customer value cycles as well as to assign customers and leads to those cycles and determining at which state within their cycle they are.

1.3.3 Telecommunications and Media

- **Network analysis:** The telecommunications industry actually already is a driver for many of the application fields described above including churn prevention, channel selection, direct marketing, customer lifetime value, and sentiment analysis. One interesting finding is that the telecommunication industry in addition can make use of another data source describing social interactions between their customers. For example, the calls between customers might describe their social connections and describing those connections in addition to their usage behavior and transactional data can easily help to improve models for churn prevention and others. If a key opinion leader decides to change to another provider, other people influenced by this person might show a higher probability to also quit their contracts as well.
- Customer service process automation: Many companies spend a lot of efforts to improve customer support and they are right to do so. If a customer is happy with the support in case of a previously bad experience, this person might be turned to a more loyal customer than before. One of the most important influence factors for the happiness with customer service is the amount of time between when a request is sent and when the answer is delivered. Text analytics can help to either answer a question with an automatically selected text block or at least by assigning the request to the right employee without any further delay. This is another example not only applicable to the telecommunications industry but to all customer-centric businesses with many customer contacts.

1.3.4 Manufacturing, Construction, and Electronics

- **Predictive maintenance:** Predictive analytics can analyze sensor data right from the production processes and machines to determine if there is an upcoming problem with a machine which is likely to lead to a problem soon. Many problems manifest themselves at early stages already with certain types of error messages or changed behavior ultimately leading to changed sensor data. Analysts can create predictive models based on the failure events in the past and the historical data before those failures. Such a model can then be used to predict if a new failure is likely to happen before the next maintenance interval and should be addressed better now than later. Those models could also deliver insights into the reasons of those failures and hence deliver a root cause analysis.
- **Patent text analysis:** Another of those examples which is actually also applicable to other industries is the analysis of patent texts. This is most often done by methods derived from text analysis and text similarities for one of two reasons: either a company would like to detect emerging trends as soon as possible or it would like to be prepared for cases in which own patents are attacked.
- Supply chain management: There are multiple application fields for predictive analytics around supply chain management. We already have discussed risk management above, for supply chain management this could mean determining which suppliers have the highest risks for failure and what would be the expected impact in case of failure. Predictive analytics can also be used for demand forecasting and hence for improving logistics but also timely negotiations with suppliers. And finally those techniques can be used to predict prices and their changes in the supply chain, again allowing for pro-active and well-informed decisions.
- **Optimizing throughput rates:** The manufacturing industry even has started to connect data mining and predictive analytics with the control centers of their factories. Based on the sensor data describing the production process itself and also the input to this process, those models find the optimal settings for the process in real time in order to optimize for quality, or higher throughput rates, or even both at the same time. In some of those cases it is important not to leave certain parameter ranges in order to not damage the involved machines and even this is possible to do with the help of advanced analytics.
- **Quality assurance:** Another application field is the prediction of the quality of the outcome of a process even before the full process has been finished. Those predictive models use the data describing the process and combine it with sensor data describing the current state of an item in order to predict the quality of the final outcome. We even have seen cases where the item was taken out of further refinement, which would just have induced additional costs for a product which will not be sold due to quality restrictions anyway. Closely connected to this are questions like anomaly detection and why a certain item has a lower quality, hence a root cause analysis.

1.4 Fundamental Terms

We managed so far to get the general idea about what are data mining and predictive analytics about. Plus we have a good feeling that those technologies are very useful for many application fields across basically all industries. In the following, we will introduce some fundamental terms which will make dealing with the problems described later in this book easier for us. You will come across these terms again and again in the RapidMiner software too, meaning it is worth becoming acquainted with the terms used even if you are already an experienced data analyst.

First of all, we can see what the two examples looked at in the previous sections, namely the toss of a coin and the falling glass, have in common. In our discussion on whether we are able to predict the end of the respective situation, we realized that knowing the influence factors as accurately as possible, such as material properties or the nature of the ground, is important. And one can even try to find an answer to the question as to whether this book will help you by recording the characteristics of yourself, the reader, and aligning them with the results of a survey of some of the past readers. These measured reader characteristics could be, for example, the educational background of the person concerned, the liking of statistics, preferences for other, possibly similar books and further features, which we could also measure as part of our survey. If we now knew such characteristics of 100 readers and had the indication as to whether you like the book or not in addition, then the further process would be almost trivial. We would also ask you the questions from our survey and measure the same features in this way and then, for example using analogy reasoning as described above, generate a reliable prediction of your personal taste. "Customers who bought this book also bought...." This probably rings a bell.

1.4.1 Attributes and Target Attributes

Whether coins, customers, or production processes, there is, as previously mentioned, the question in all scenarios as to the characteristics or features of the respective situation. We will always speak of **attributes** in the following when we mean such describing factors of a scenario. This is also the term that is always used in the RapidMiner software when such describing features arise. There are many synonyms for this term and depending on your own background you might have already come across different terms instead of "attribute", for example

- characteristic,
- feature,
- influence factor (or just factor),
- indicator,
- variable, or
- signal.

We have seen that description by attributes is possible for processes and also for situations. This is necessary for the description of technical processes for example and the thought of the falling glass is not too far off here. If it is possible to predict the outcome of such a situation, then why not also the quality of a produced component? Or the imminent

failure of a machine? Other processes or situations which have no technical reference can also be described in the same way. How can I predict the success of a sales or marketing promotion? Which article will a customer buy next? How many more accidents will an insurance company probably have to cover for a particular customer or customer group?

We shall use such a customer scenario in order to introduce the remaining important terms. Firstly, because humans are famously better at understanding examples about other humans and secondly, because each enterprise probably has information, i.e., attributes, regarding their customers and most readers can therefore relate to the examples immediately. The attributes available as a minimum, which just about every enterprise keeps about its customers, are for example address data and information as to which products or services the customer has already purchased. You would be surprised what forecasts can be made even from such a small number of attributes.

Let us look at an (admittedly somewhat contrived) example. Let us assume that you work in an enterprise that would like to offer its customers products in the future which are better tailored to their needs. Within a customer study of only 100 of your customers some needs became clear, which 62 of these 100 customers share all the same. Your research and development department got straight to work and developed a new product within the shortest time, which would satisfy these new needs better. Most of the 62 customers with the relevant needs profile are impressed by the prototype in any case, although most of the remaining participants of the study only show a small interest as expected. Still, a total of 54 of the 100 customers in the study said that they found the new product useful. The prototype is therefore evaluated as successful and goes into production—now only the question remains as to how, from your existing customers or even from other potential customers, you are going to pick out exactly the customers with whom the subsequent marketing and sales efforts promise the greatest success. You would therefore like to optimize your efficiency in this area, which means in particular ruling out such efforts from the beginning which are unlikely to lead to a purchase. But how can that be done? The need for alternative solutions and thus the interest in the new product arose within the customer study on a subset of your customers. Performing this study for all your customers is much too costly and so this option is closed to you. And this is exactly where data mining can help. Let us first look at a possible selection of attributes regarding your customers:

- Name
- Address
- Sector
- Subsector
- Number of employees
- Number of purchases in product group 1
- Number of purchases in product group 2

The number of purchases in the different product groups means the transactions in your product groups which you have already made with this customer in the past. There can of course be more or less or even entirely different attributes in your case, but this is irrelevant at this stage. Let us assume that you have the information available regarding these attributes for every one of your customers. Then there is another attribute which we can look at for our concrete scenario: Whether the customer likes the prototype or not. This attribute is of course only available for the 100 customers from the study; the information on

this attribute is simply unknown for the others. Nevertheless, we also include the attribute in the list of our attributes:

- Prototype positively received?
- Name
- Address
- Sector
- Subsector
- Number of employees
- Number of purchases in product group 1
- Number of purchases in product group 2

If we assume you have thousands of customers in total, then you can only indicate whether 100 of these evaluated the prototype positively or not. You do not yet know what the others think, but you would like to! The attribute "prototype positively received" thus adopts a special role, since it identifies every one of your customers in relation to the current question. We therefore also call this special attribute a **label**, since it sticks to your customers and identifies them like a brand label on a shirt or even a note on a pin board. You will also find attributes which adopt this special role in RapidMiner under the name "label". The goal of our efforts is to fill out this particular attribute for the total quantity of all customers. We might therefore also speak of the **target attribute** instead of the term "label" since our target is to create a model which predicts this special attribute from the values of all others. You will also frequently discover the term **target variable** in the literature, which means the same thing.

1.4.2 Concepts and Examples

The structuring of your customers' characteristics by attributes, introduced above, already helps us to tackle the problem a bit more analytically. In this way we ensured that every one of your customers is represented in the same way. In a certain sense we defined the type or **concept** "customer", which differs considerably from other concepts such as "falling articles" in that customers will typically have no material properties and falling articles will only rarely buy in product group 1. It is important that, for each of the problems in this book (or even those in your own practice), you first define which concepts you are actually dealing with and which attributes these are defined by.

We implicitly defined above, by indicating the attributes name, address, sector, etc., and in particular the purchase transactions in the individual product groups, that objects of the concept "customer" are described by these attributes. Yet this concept has remained relatively abstract so far and no life has been injected into it yet. Although we now know in what way we can describe customers, we have not yet performed this for specific customers. Let us look at the attributes of the following customer, for example:

- Prototype positively received: yes
- Name: Miller Systems Inc.
- Address: 2210 Massachusetts Avenue, 02140 Cambridge, MA, USA

- Sector: Manufacturing
- Subsector: Pipe bending machines
- Number of employees: >1000
- Number of purchases in product group 1: 5
- Number of purchases in product group 2: 0

We say that this specific customer is an **example** for our concept "customer". Each example can be characterized by its attributes and has concrete values for these attributes which can be compared with those of other examples. In the case described above, Miller Systems Inc. is also an example of a customer who participated in our study. There is therefore a value available for our target attribute "prototype positively received?" Miller Systems Inc. was happy and has "yes" as an attribute value here, thus we also speak of a **positive example**. Logically, there are also negative examples and examples which do not allow us to make any statement about the target attribute.

1.4.3 Attribute Roles

We have now already become acquainted with two different kinds of attributes, i.e., those which simply describe the examples and those which identify the examples separately. Attributes can thus adopt different roles. We have already introduced the role "label" for attributes which identify the examples in any way and which must be predicted for new examples that are not yet characterized in such a manner. In our scenario described above, the label still describes (if present) the characteristic of whether the prototype was received positively.

Likewise, there are for example roles, the associated attribute of which serves for clearly identifying the example concerned. In this case the attribute adopts the role of an identifier and is called **ID** for short. You will find such attributes named ID in the RapidMiner software also. In our customer scenario, the attribute "name" could adopt the role of such an identifier.

There are even more roles, such as those with an attribute that designates the weight of the example with regard to the label. In this case the role has the name **Weight**. Attributes without a special role, i.e., those which simply describe the examples, are also called **regular** attributes and just leave out the role designation in most cases. Apart from that you have the option in RapidMiner of allocating your own roles and of therefore identifying your attributes separately in their meaning. Please note that for most data mining tasks in RapidMiner the regular attributes are used as input to the method, for example, to create a model predicting the attribute with role label.

1.4.4 Value Types

As well as the different roles of an attribute, there is also a second characteristic of attributes which is worth looking at more closely. The example of Miller Systems Inc. above defined the respective values for the different attributes, for example "Miller Systems Inc." for the attribute "Name" and the value "5" for the number of past purchases in product group 1. Regarding the attribute "Name", the concrete value for this example is therefore arbitrary free text to a certain extent; for the attribute "number of purchases in product group 1" on the other hand, the indication of a number must correspond. We call the indication whether the values of an attribute must be in text or numbers the Value Type of an attribute.

14

In later chapters we will become acquainted with many different value types and see how these can also be transformed into other types. For the moment we just need to know that there are different value types for attributes and that we speak of value type **text** in the case of free text, of the value type **numerical** in the case of numbers and of the value type **nominal** in the case of only few values being possible (like with the two possibilities "yes" and "no" for the target attribute).

Please note that in the above example the number of employees, although really of numerical type, would rather be defined as nominal, since a size class, i.e., ">1000" was used instead of an exact indication like 1250 employees.

1.4.5 Data and Meta Data

We want to summarize our initial situation one more time. We have a Concept "customer" available, which we will describe with a set of attributes:

- Prototype positively received? [Label; Nominal]
- Name [Text]
- Address [Text]
- Sector [Nominal]
- Subsector [Nominal]
- Number of employees [Nominal]
- Number of purchases in product group 1 [Numerical]
- Number of purchases in product group 2 [Numerical]

The attribute "Prototype positively received?" has a special **Role** among the attributes; it is our **Target Attribute** or **Label** here. The target attribute has the value type **Nominal**, which means that only relatively few characteristics (in this case "yes" and "no") can be accepted. Strictly speaking it is even binominal, since only two different characteristics are permitted. The remaining attributes all have no special role, i.e., they are **regular** and have either the value type **Numerical** or **Text**. The following definition is very important, since it plays a crucial role in a successful professional data analysis: This volume of information, which describes a concept, is also called meta data, since it represents data about the actual data.

Our fictitious enterprise has a number of Examples for our concept "customer", i.e., the information which the enterprise has stored for the individual attributes in its customer database. The goal is now to generate a prediction instruction from the examples for which information is available concerning the target attribute, which predicts for us whether the remaining customers would be more likely to receive the prototype positively or reject it. The search for such a prediction instruction is one of the tasks which can be performed with data mining.

However, it is important here that the information for the attributes of the individual examples is in an ordered form, so that the data mining method can access it by means of a computer. What would be more obvious here than a table? Each of the attributes defines a column and each example with the different attribute values corresponds to a row of this table. For our scenario this could look like Table 1.1 for example.

We call such a table an **Example Set**, since this table contains the data for all the attributes of our examples. In the following, and also within RapidMiner, we will use the

Prototype positively received?	Name	Address	Sector	Sub- sector	Number of em- ployees	No of pur- chases group 1	No of pur- chases group 2	
Yes	Miller Systems Inc.	2210 Mas- sachusetts Ave, Cam- bridge, MA, USA	Manufacturing	Pipe bending ma- chines	>1000	5	0	
?	Smith Paper	101 Hunt- ington Ave, Boston, MA, USA	IT	Supplies	600-1000	3	7	
No	Meyer Inc.	1400Com-merceSt,Dallas,TX,USA	Retail	Textiles	<100	1	11	

TABLE 1.1: A set of examples with values for all attributes together with label values.

terms **Data**, **Dataset**, and **Example Set** synonymously. A table with the appropriate entries for the attribute values of the current examples is always meant in this case. It is also such data tables which have taken their name from data analysis or data mining. Note:

Data describes the objects of a concept, **Meta Data** describes the characteristics of a concept (and therefore also of the data).

Most data mining methods expect the examples to be given in such an attribute value table. Fortunately, this is the case here and we can spare ourselves any further data transformations. In practice, however, this is completely different and the majority of work during a data analysis is time spent transferring the data into a format suitable for data mining. These transformations are therefore dealt with in detail in later chapters when we are discussing the different use cases.

1.4.6 Modeling

Once we have the data regarding our customers available in a well-structured format, we can then finally replace the unknown values of our target attribute with the prediction of the most probable value by means of a data mining method. We have numerous methods available here, many of which, just like the analogy reasoning described at the beginning or the generating of rules of thumb, are based on human behavior. We call the use of a data mining method "to model" and the result of such a method, i.e., the prediction instruction, is a model. Just as data mining can be used for different issues, this also applies for models. They can be easy to understand and explain the underlying processes in a simple manner. Or they can be good to use for prediction in the case of unknown situations. Sometimes both apply, such as with the following model for example, which a data mining method could have supplied for our scenario:

If the customer comes from urban areas, has more than 500 employees and if at least 3 purchases were transacted in product group 1, then the probability of this customer being interested in the new product is high.

Such a model can be easily understood and may provide a deeper insight into the underlying data and decision processes of your customers. And in addition, it is an **operational model**, i.e., a model which can be used directly for making a prediction for further customers. The company "Smith Paper" for example satisfies the conditions of the rule above and is therefore bound to be interested in the new product—at least there is a high probability of this. Your goal would therefore have been reached and by using data mining you would have generated a model which you could use for increasing your marketing efficiency: Instead of just contacting all existing customers and other candidates without looking, you could now concentrate your marketing efforts on promising customers and would therefore have a substantially higher success rate with less time and effort. Or you could even go a step further and analyze which sales channels would probably produce the best results and for which customers.

In the rest of this book we will focus on further uses of data mining and at the same time practice transferring concepts such as customers, business processes, or products into attributes, examples, and datasets. This will train the eye to detect further possibilities of application tremendously and will make analyst life much easier for you later on. In the next chapter though, we would like to spend a little time on RapidMiner and give a small introduction to its use, so that you can implement the following examples immediately.

Chapter 2

Getting Used to RapidMiner

Ingo Mierswa

Rapid-I, Dortmund, Germany

2.1	Introd	luction	19	
2.2	First S	t Start		
2.3	Design	1 Perspective	21	
2.4	Buildi	ng a First Process	23	
	2.4.1	Loading Data	24	
	2.4.2	Creating a Predictive Model	25	
	2.4.3	Executing a Process	28	
	2.4.4	Looking at Results	29	

2.1 Introduction

We will use the open source data mining solution RapidMiner for the practical exercises in this book and for demonstrating the use cases for data mining, predictive analytics, and text mining discussed here. Before we can work with RapidMiner, you of course need to download and install the software first. You will find RapidMiner in the download area of the Rapid-I website: http://www.rapid-i.com. If you want to follow the use cases in this book, it is highly recommended to download the appropriate installation package for your operating system and install RapidMiner according to the instructions on the website now. All usual Windows versions are supported as well as Macintosh, Linux or Unix systems. Please note that an up-to-date Java Runtime Environment (JRE) with at least version 7 is needed for all non-Windows systems. You can find such a JRE for example at http://www.java.com/.

2.2 First Start

If you are starting RapidMiner for the first time, you will be welcomed by the so-called Welcome Perspective. The lower section shows current news about RapidMiner, if you have an Internet connection. The list in the center shows the analysis processes recently worked on. This is useful if you wish to continue working on or execute one of these processes. You can open a process from this list to work on or execute it simply by double clicking. The upper section shows typical actions which you as an analyst will perform frequently after starting RapidMiner. Here are the details of these:

1. New Process: Starts a new analysis process. This will be the most often used selection

for you in the future. After selecting this, RapidMiner will automatically switch to the Design perspective (explained below).

- 2. Open Recent Process: Opens the process which is selected in the list below the actions. Alternatively, you can open this process by double-clicking inside the list. Either way, RapidMiner will then automatically switch to the Design Perspective.
- 3. Open Process: Opens the repository browser and allows you to select a process to be opened within the process Design Perspective.
- 4. Open Template: Shows a selection of different pre-defined analysis processes, which can be configured in a few clicks.
- 5. Online Tutorial: Starts a tutorial which can be used directly within RapidMiner and gives an introduction to some data mining concepts using a selection of analysis processes. Recommended if you have a basic knowledge of data mining and are already familiar with the fundamental operation of RapidMiner.



FIGURE 2.1: Welcome Perspective of RapidMiner.

At the right-hand side of the toolbar inside the upper section of RapidMiner, you will find three icons which switch between the individual RapidMiner perspectives. A **perspective** consists of a freely configurable selection of individual user interface elements, the so-called views. These can also be arranged however you like. In the Welcome Perspective there is only one view; one preset at least, namely the welcome screen, which you are looking at now. You can activate further views by accessing the "View" menu. Sometimes you may inadvertently delete a view or the perspective is unintentionally moved into particularly unfavorable positions. In this case the "View" menu can help, because apart from the possibility of reopening closed views via "Show View", the original state can also be recovered at any time via "Restore Default Perspective".

2.3 Design Perspective

You will find an icon for each (pre-defined) perspective within the right-hand area of the toolbar:



FIGURE 2.2: Toolbar icons for perspectives.

The icons shown here take you to the following perspectives:

- 1. *Design Perspective:* This is the central RapidMiner perspective where all analysis processes are created and managed.
- 2. *Result Perspective:* If a process supplies results in the form of data, models, or the like, then RapidMiner takes you to this Result Perspective, where you can look at several results at the same time as normal thanks to the views.
- 3. *Welcome Perspective:* The Welcome Perspective already described above, in which RapidMiner welcomes you with after starting the program.

You can switch to the desired perspective by clicking inside the toolbar or alternatively via the menu entry "View" – "Perspectives" followed by the selection of the target perspective. RapidMiner will eventually also ask you automatically if switching to another perspective seems a good idea, e.g., to the Result Perspective on completing an analysis process.

Now switch to the Design Perspective by clicking in the toolbar. This is the major working place for us while using RapidMiner. Since the Design Perspective is the central working environment of RapidMiner, we will discuss all parts of the Design Perspective separately in the following and discuss the fundamental functionalities of the associated views. In any case, you should now see the screen on the next page:

All work steps or building blocks for different data transformation or analysis tasks are called **operators** in RapidMiner. Those operators are presented in groups in the Operator View on the left side. You can navigate within the groups in a simple manner and browse in the operators provided to your heart's desire. If RapidMiner has been extended with one of the available extensions, then the additional operators can also be found here. Without extensions you will find at least the following groups of operators in the tree structure:

- *Process Control:* Operators such as loops or conditional branches which can control the process flow.
- *Utility:* Auxiliary operators which, alongside the operator "Subprocess" for grouping subprocesses, also contain the important macro-operators as well as the operators for logging.
- *Repository Access:* Contains the two operators for read and write access in repositories.
- *Import:* Contains a large number of operators in order to read data and objects from external formats such as files, databases, etc.



FIGURE 2.3: Design Perspective of RapidMiner.

- *Export:* Contains a large number of operators for writing data and objects into external formats such as files, databases, etc.
- Data Transformation: Probably the most important group in the analysis in terms of size and relevance. All operators are located here for transforming both data and meta data.
- *Modeling:* Contains the actual data mining process, such as classification methods, regression methods, clustering, weightings, methods for association rules, correlation and similarity analyses as well as operators, in order to apply the generated models to new datasets.
- *Evaluation:* Operators using which one can compute the quality of a modeling and thus for new data, e.g., cross-validations, bootstrapping, etc.

You can select operators within the Operators View and add them in the desired place in the process by simply dragging them from the Operators View and dropping them into the large white area in the center, the so-called Process View. Every analysis in RapidMiner is a **process**, and every process consists of one or several steps which are the operators.

Depending on your settings, those new operators might be directly connected with existing operators as suitably as possible on the basis of the available meta data information. If this is not happening or the automatically inserted connection is not desired, you can delete the connection by selecting them and pressing the Delete key or by pressing the Alt key while clicking on any of the connection **ports**. Ports are the round bubbles on the sides of the operators and they are used to define the data flow through your analytical processes. You can insert new connections by either clicking on the source port and then clicking again on a target port or by dragging a line between those ports.

Later on, when you have successfully defined your first RapidMiner processes, a typical result might look like in the image on the following page:

You could now simply try and select a few operators from the Operator View and drag them into the Process View. Connect their ports, even if this is probably not leading

22



FIGURE 2.4: A typical process in RapidMiner consists of several operators. The data flow is defined by the connected ports, for example, the "Read Excel" operator loads a dataset from an Excel file and delivers it to the succeeding "Select Attributes" operator.

to working processes, and get familiar with the user interface of RapidMiner. In order to edit parameters you must select an individual operator. You will recognize the operator currently selected by its orange frame as well as its shadow. If you wish to perform an action for several operators at the same time, for example moving or deleting, please select the relevant operators by dragging a frame around these. In order to add individual operators to the current selection or exclude individual operators from the current selection, please hold the CTRL key down while you click on the relevant operators or add further operators by dragging a frame. You can also move operators around by selecting them and dragging them in the Process View. You will notice that the parameters in the **Parameter View** on the right side of RapidMiner changes sometimes if you select different operators. As you can see, most operators provide a set of parameters which control the actual working mode of the respective operator. You will find much information about RapidMiner and its user interface in the RapidMiner User Manual available at

http://docs.rapid-i.com

2.4 Building a First Process

One of the first steps in a process for data analysis is usually to load some data into the system. RapidMiner supports multiple methods for accessing datasets. It supports more than 40 different file types and of course all major database systems. If the data is not originally stored in a relational database system, the best approach is to import the data first into the RapidMiner repository. Please follow the instructions from the RapidMiner manual for more information or just try to import a dataset, for example an Excel file, with "File" – "Import Data". Later you will also realize that there are dozens of operators for data import in the operator group "Import", which can also be used directly as part of the process.

2.4.1 Loading Data

In the following we assume that you have managed to import the data into the Rapid-Miner repository and hence we will retrieve the data from there. If you are loading the data from a database or file, your following steps are at least similar to those described below.

It is always recommended to use the repository whenever this is possible instead of files. This will allow RapidMiner to get access to the describing meta data and will ease process design a lot. We will now create the beginning of a data analysis process and will add a first data mining technique using this data. The very first operation in our process should be to load the data from the repository again in order to make it available for the next analysis steps:

1. Go to the Repositories view and open the repository **Samples** delivered with Rapid-Miner. Click on the small plus sign in front of this repository. You should now see two folders named **data** and **processes**. Open the data folder and you will find a collection of datasets coming together with RapidMiner. Click on the dataset named **Iris** and drag it onto the large white view named **Process** in the center of your frame. After releasing the dataset somewhere on the white area, it should be transformed into an operator named **Retrieve** with a bluish output port on the right side. RapidMiner automatically has transformed the dataset into an operator loading the dataset. If you click on the operator, you can see a parameter in the **Parameters** view pointing to the data location. The Retrieve operator in general, well, retrieves objects from a repository and makes them available in your process.



FIGURE 2.5: Drag the Iris dataset into the process view in order to create a new operator loading this dataset during process execution.

Background Information – Iris Data:

You might already be familiar with the "Iris" dataset since it is a well-known dataset among many data analysts. If not, here is the basic idea: the dataset describes 150 Iris plants with four attributes: sepal-length, sepal-width, petal-length, and petal-width. Since only experts for Iris plants actually understand the meaning of those attributes we will refer to those attributes with "a1" to "a4". And there is a fifth column describing to which class of Irises each of the 150 plants belong. There are three options: Iris-setosa, Iris-versicolor, and Iris-virginica. Each of those three classes is represented by 50 plants in the dataset. The goal is now to find a classification model using the measured lengths and widths (a1 to a4) to predict the class of the plant. This would
allow the classification of those plants even by non-experts like myself.

2. Maybe the output was automatically connected to one of the result ports named **res** on the right side of the Process area. If that is the case, you are done already. If not—and this depends on your current program settings—click on the output port of the Retrieve operator and then click on the first res port on the right side. Alternatively, you can also drag a connection line between the two ports. Your process should now look like the following:



FIGURE 2.6: The probably most simple process which can be created with RapidMiner: It just retrieves data from a repository and delivers it as a result to the user to allow for inspection.

All results which are delivered at a result port named **res** will be delivered as a result to the user (or other systems if the process execution is integrated into other systems). The next step would be to create a decision tree on the Iris dataset and also deliver this model to the user.

2.4.2 Creating a Predictive Model

We have seen above how we create a new process just loading a dataset. The next step will be to create a predictive model using this dataset. This model predicts a categorical or nominal value, hence we could also say the model should describe rules which allow us to assign one of the three classes to new and unseen data describing new plants. We refer to this type of modeling as **classification**.

Adding a modeling technique to your process so that it calculated a predictive model is actually very easy. Just follow the following steps for creating such a model:

Go to the Operators view and open the operator group **Modeling**, **Classification**, and **Regression**, and then **Tree Induction**. You should now see an operator named **Decision Tree**. Click on it and drag it to your process, somewhere to the right of your initial retrieve operator.

You now only have to create the necessary connections. The dataset should be delivered to the modeling operator which is delivering a model then. However, you can also deliver the dataset itself to the user if you also connect the data port with one of the result ports. The complete process should look like the following figure.

In the next section we will learn how to execute this process and how to inspect the created results.



FIGURE 2.7: Drag the operator named "Decision Tree" into your process.



FIGURE 2.8: The complete process consisting of data loading and model creation.

Background Information: Decision Trees

Decision trees are probably one of the most widely used techniques in data mining. One of the biggest advantages is that they are easy to use and, maybe even more important, also easy to understand even by non-experts. The basic idea behind decision trees is a so-called divide-and-conquer approach. In each step the dataset is divided into different parts while each part should better represent one of the possible classes. The final result will be a treestructure where each inner node represents a test for the value of a particular attribute and each leaf is representing the decision for a particular class. A new and unknown case is then routed down the tree until it reaches one of the leaves.

For each node we have two options depending on the value type of the attribute used at this node. For nominal attributes, the number of children is usually equal to the number of possible values for this attribute. If we are using a nominal attribute for a test in one of the inner nodes, this means that the dataset is at this stage basically divided according to the different values of this attribute. Hence, a nominal attribute will not get tested more than once since all examples further down the tree will have the same value for this particular attribute. This is different for numerical attributes: Here, we usually test if the attribute value is greater or less than a determined constant. The attribute may get tested several times for different constants.

The strategy for the construction of a decision tree is top-down in a recursive divide-andconquer fashion. The first step is to select one of the attributes for the root node. Then we create a branch for each possible attribute value and split instances into subsets according to the possible values, i.e., we will get one subset for each branch extending from the node. Finally we will repeat these steps recursively for each branch but only use instances that reach the branch now. We can stop this process if all instances have the same class.

The major question is: How should we select which attribute should be tested next? Our goal is to get the smallest tree possible since we believe that a small tree manages to explain the data best and in a general fashion leading to fewer errors for unseen data compared to a more complex tree. There is not a strict optimal solution for this task but a widely used heuristic: we should choose the attribute that produces the "purest" subsets of data with respect to the label attribute. A very popular choice for this type of measurement is the so-called **information gain**. The information gain basically increases with the average purity of the subsets and hence we should choose the attribute which gives us the greatest information gain at each level.

We will now discuss how we can measure information at all in order to be able to calculate which decision delivers the highest gain in information. We will see now that information can be easily measured in bits and that a concept exists for calculating this amount: the **entropy**. The entropy measures the information required in bits (can also mean fractions of bits) to predict the outcome of an event if the probability distribution of this event is given. The formula is as follows:

Entropy $(p_1,\ldots,p_n) = -p_1 \ ld \ (p_1) - \ldots - p_n \ ld \ (p_n)$

 p_1 to p_n are the possible probabilities for all outcomes while ld is the logarithm dualis, i.e., the logarithm with base 2. The formula is easy to understand with a small example. Let's assume we are tossing an unbiased coin where each side can be shown with a probability of 0.5. In this case the entropy is -0.5 * ld(0.5) - 0.5 * ld(0.5) = 1. The necessary information to decide which side of a tossed coin is actually shown is hence 1 bit, which is the perfect amount for a binary decision with only two possible outcomes which are equally likely.

Now let's assume that the coin is biased and shows "head" with a probability of 75% and "tails" with 25%. In this case the entropy would be calculated as -0.75 * ld(0.75) - 0.75 *

ld(0.75) = 0.81 and hence we need less information to decide. This is only natural since we already expected "head" with a higher probability.

Algorithms for creating decision trees make use of the notion of entropy for selecting the optimal attribute for each split, i.e., inner node. The entropy for the class distribution is first calculated for each possible value of a nominal attribute. Let's do this for an example and assume that we have an attribute with two possible values A and B. We have 10 examples with value A and 20 examples with value B. For those with value A we have 3 examples with class X and 7 examples with class Y. And for those with value B we have 15 examples with class X and 5 examples with class Y. We can now calculate the entropy values based on those class distributions in the subsets defined by those examples having values A and B respectively for this attribute:

- 1. Entropy $A \rightarrow$ entropy of class distribution (3/10, 7/10) = 0.881 bits
- 2. Entropy $B \rightarrow$ entropy of class distribution (16/20, 4/20) = 0.322 bits

For numerical attributes, each possible threshold for a numerical split is tested instead. In both cases the total information needed given a specific attribute is the weighted average of all entropy values; in this case this would be

(10 * Entropy A + 20 * Entropy B) / 30 = 0.508.

In order to calculate the **information gain** by this split, we have to compare the entropy of the class distribution before the split with the entropy after the split, i.e.,in our example:

gain = info (19/30, 11/30) - weighted entropy = 0.948 - 0.508 = 0.44

The information gain for performing this split hence is 0.44. We can now calculate the gain for all possible attributes and compare all gains to find the one delivering the highest value. This attribute would then be used for the next split. The dataset is divided accordingly and the complete process is repeated recursively until all subsets only contain examples for a single class.

The section above has explained the basic ideas behind decision trees, but be assured that there is more about decision trees beyond this basic approach. Instead of the information gain, one could use different measurements including the so-called **gain ratio** which prefers splits with less possible branches. Another important addition is **pruning** the tree after the growing phase in order to make it more robust for unseen data points. Please refer to the standard literature about data mining for getting more information about additions to the basic decision tree algorithm.

2.4.3 Executing a Process

Now we are ready and want to execute the process we have just created for the first time. The status indicators of all used operators should now be yellow or green (the small traffic lights in each operator box) in the best case, and there should be no entries in the Problems View before you start executing a process. In such a case it should be possible to execute our process currently consisting of only one operator without any problems. However, the problems in the Problems view and also the red traffic lights only indicate that there might be a potential problem—you might be able to execute a process even if RapidMiner detects a potential problem. Just execute the process and check if it works despite the complaint as described below.

You have the following options for starting the process:

- 1. Press the large play button in the toolbar of RapidMiner.
- 2. Select the menu entry "Process" \rightarrow "Run".
- 3. Press F11.



FIGURE 2.9: Press the Play icon in order to execute your process.

While a process is running, the status indicator of the operator being executed in each case transforms into a small green play icon. This way you can see what point the process is currently at. After an operator has been successfully executed, the status indicator then changes and stays green—until for example you change a parameter for this operator: Then the status indicator will be yellow. The same applies for all operators that follow. This means you can see very quickly which operators a change could have an effect on.

The process defined above only has a short runtime and so you will hardly have the opportunity to pause the running process. In principle, however, you can briefly stop a running process with the pause symbol, e.g., in order to see an intermediate result. The operator currently being executed is then finished and the process is then stopped. You can recognize a process that is still running but currently paused by the fact that the color of the play icon changes from blue to green.

Press the play button again to continue executing the process further. If you do not wish to merely pause the process but to abort it completely, then you can press the stop button. Just like when pausing, the operator currently being executed is finished and the process fully aborted immediately after. Please note that you can switch to the Design Perspective immediately after aborting the process and make changes to processes, even if the execution of the current operator is being finished in the background. You can even start further processes and do not need to wait for the first process to be completed.

Note: The operator being executed is always completed if you abort. This is necessary to ensure a sound execution of operators. However, completing an operator may need much more time in individual cases and also require other resources such as storage space. So if when aborting very complex operators you can see this taking hours and requiring additional resources, then your only option is to restart the application.

2.4.4 Looking at Results

After the process was terminated, RapidMiner should automatically have switched to the Result Perspective. If this was not the case, then you probably did not connect the output port of the last operator with one of the result ports of the process on the right-hand side. Check this and also check for other possible errors, taking the notes in the Problems View into consideration.

Feel free to spend a little time with the results. The process above should have delivered a dataset and a decision tree used to predict the label of the dataset based on the attributes' values. You can inspect the data itself as well as the meta data of this dataset and try out some of the visualizations in the plot view. You can also inspect the decision tree and try to understand if this makes sense to you. If you wish to return to the Design Perspective, then you can do this at any time using the switching icons at the right of the toolbar.

Tip: After some time you will want to switch frequently between the Design Perspective

and the Result Perspective. Instead of using the icon or the menu entries, you can also use keyboard commands F8 to switch to the Design Perspective and F9 to switch to the Result Perspective.

What does that result mean to us? We now have managed to load a dataset from the RapidMiner repository and then we have built the first predictive model based on this data. Furthermore, we got a first feeling about how to build RapidMiner processes. You are now ready to learn more about the use cases for data mining and how to build corresponding processes with RapidMiner. Each of the following chapters will describe a use case together with the data which should be analyzed. At the same time each chapter will introduce new RapidMiner operators to you which are necessary to successfully solve the tasks at hand.



FIGURE 2.10: The decision tree described the rules how to assign the different classes to a new plant. First we check the value of a3 and if it small enough, we assign the blue class Iris-setosa to this plant. If not, we additionally check the value for a4 and if this is high we assign the red class Iris-virginica to it. Finally, we will check for attribute a3 again and decide for Iris-versicolor if the value is small enough.

Part II

Basic Classification Use Cases for Credit Approval and in Education

Chapter 6

Naïve Bayes Classificaton II

M. Fareed Akhtar

Fastonish, Australia

6.1	Datas	et	65
	6.1.1	Nursery Dataset	65
	6.1.2	Basic Information	65
	6.1.3	Examples	65
	6.1.4	Attributes	66
6.2	Opera	tors in this Use Case	67
	6.2.1	Read Excel Operator	67
	6.2.2	Select Attributes Operator	67
6.3	Use C	ase	67
	6.3.1	Data Import	67
	6.3.2	Pre-processing	69
	6.3.3	Model Training, Testing, and Performance Evaluation	69
	6.3.4	A Deeper Look into the Naïve Bayes Algorithm	71

The use case of this chapter applies the Naïve Bayes operator on the Nursery dataset (overview of the Naïve Bayes algorithm has been discussed in the previous chapter). The working of the Naïve Bayes operator is also discussed in detail. The purpose of this use case is to rank applications for nursery schools using the Naïve Bayes algorithm. The operators explained in this chapter are: Read Excel and Select Attributes operators.

6.1 Dataset

6.1.1 Nursery Dataset

This dataset has been taken from the UCI repositories. This dataset can be accessed through this link: http://archive.ics.uci.edu/ml/datasets/Nursery.

6.1.2 Basic Information

Nursery Database was derived from a hierarchical decision model originally developed to rank applications for nursery schools. It was used during several years in the 1980s when there was excessive enrollment in these schools in Ljubljana, Slovenia, and the rejected applications frequently needed an objective explanation.

6.1.3 Examples

This dataset has 12,960 examples and there are no missing values in this dataset.

6.1.4 Attributes

This dataset has 9 attributes (including the label attribute). The dataset comes with some basic information, but the type and role of attributes is set by the user of the dataset. Even the attribute names are specified by the user. Here is an explanation of the attributes of this dataset:

- 1. **Parents:** This attribute has information about the parents of the child. It has three possible values: usual, pretentious, and great_pret. As this attribute can have three possible values, the type of this attribute should be set to polynominal in RapidMiner.
- 2. **Has_nur:** This attribute has information about the nursery of the child. It has five possible values: proper, less_proper, improper, critical, and very_crit. As this attribute can have five possible values, the type of this attribute should be set to polynominal in RapidMiner.
- 3. Form: This attribute has information about the form filled out by the applicants. It has four possible values: complete, completed, incomplete, and foster. As this attribute can have four possible values, the type of this attribute should be set to polynomial in RapidMiner.
- 4. Children: This attribute has information about the number of children of the applicant. It has four possible values: {1, 2, 3, more}. As this attribute can have four possible values, the type of this attribute should be set to polynominal in RapidMiner.
- 5. **Housing**: This attribute has information about the housing standard of the applicant. It has three possible values: convenient, less_conv, and critical. As this attribute can have three possible values, the type of this attribute should be set to polynominal in RapidMiner.
- 6. **Finance:** This attribute has information about the financial standing of the applicant. It has two possible values: convenient, inconv. As this attribute can have only two possible values, the type of this attribute should be set to binominal in RapidMiner.
- 7. Social: This attribute has information about the social structure of the family. It has three possible values: nonprob, slightly_prob, and problematic. As this attribute can have three possible values, the type of this attribute should be set to polynominal in RapidMiner.
- 8. **Health:** This attribute has information about the health picture of the family. It has three possible values: recommended, priority, and not_recom. As this attribute can have three possible values, the type of this attribute should be set to polynominal in RapidMiner.
- 9. Rank (label attribute): This attribute specifies the rank of the application. It has five possible values: not_recom, recommend, very_recom, priority, and spec_prior. As this attribute has five possible values its type should be set to polynominal in RapidMiner. The role of this attribute should be set to label because this is the target attribute or the attribute whose value will be predicted by the classification algorithms. The role of all other attributes should be set to regular.

More information about this dataset can be obtained from UCI repositories.

6.2 Operators in this Use Case

6.2.1 Read Excel Operator

For applying any algorithm on a dataset the first step is to import the dataset. Importing the dataset means loading the dataset into RapidMiner. RapidMiner provides numerous operators for loading datasets. These operators can be found at the 'Import/Data' section in the Operators window. Data is available in different formats. Therefore, RapidMiner provides different operators for importing data in these different formats.

Mostly data is stored in CSV files, Excel files, or in databases. To access such datasets, RapidMiner provides operators like Read CSV, Read Excel, and Read Database.

The Read Excel operator imports data from Microsoft Excel spreadsheets. The user has to define which of the spreadsheets in the workbook should be used as data table. The table must have a format so that each row is an example and each column represents an attribute. The data table can be placed anywhere on the sheet.

The easiest and shortest way to import an Excel file is to use *import configuration wizard* from the *Parameters* View. The import configuration wizard, Excel file, first row as names, and dataset meta data information parameters of the Read Excel operator are very similar to corresponding parameters in the Read CSV operator. Other important parameters of the Read Excel operator include the *sheet number* parameter which specifies the number of the sheet which contains the required dataset. The *imported cell range* parameter specifies the range of cells to be imported. It is specified in "xm:yn" format where "x" is the column of the first cell of range, "m" is the row of the first cell of range, "y" is the column of the last cell of range, "n" is the row of the last cell of range. For example, "A1:E10" will select all cells of the first five columns from row 1 to 10.

6.2.2 Select Attributes Operator

This is a very powerful operator for selecting required attributes of the given dataset because it gives many different options for selecting attributes. Attributes can be selected by simply mentioning their names. Attributes of a particular type or block type can also be selected. Attributes can be selected on the basis of missing values or even attribute value filters. The most powerful option for selecting attributes is through regular expression. Moreover, exceptions can be provided for any method of attribute selection.

6.3 Use Case

The main goal of this process is to have a deeper look at the working of the Naïve Bayes algorithm. The Read Excel operator is used for importing the Nursery dataset. The Split Data operator is used for splitting the data into training and testing datasets. Finally, the Performance (Classification) operator is used for evaluating the performance of the model.



FIGURE 6.1: Workflow of the process.

6.3.1 Data Import

The Read Excel operator is used for importing data in this process. The data at the following url is stored in an Excel file. http://archive.ics.uci.edu/ml/machine-learning-databases/nursery/nursery.data.

This Excel file is then used as the data source in RapidMiner. The Read Excel operator is used for loading data from an Excel file. The *Import Configuration Wizard* is used for loading the data in this process. The steps of this configuration wizard are explained below.

Step 1: Select Excel file.

🚯 Data import wizard -	Step 1 of 4			×	🤳 R	ead Excel
This wize	rd quides you to import your da	ita			🏸 Import Confi	guration Wizard
Step 1: P	lease select the file that should	be imported.			excel file	WWWursery.xls
Look In: 🔒 N		•	🔹 🍫 🎓 😂	Q -	sheet number	1
Bookmarks	File Name	Size	Type Last	Modified	imported cell range	A1:I12960
🚖 Last Directory	Nursery.xls	1 MB N	Aicrosoft Offic Mar 3	1, 2012	encoding	SYSTEM -
					🗌 first row as name	s
					annotations	Edit List (0)
File Name: Nursery	xls				date format	•
Files of Type: Excel F	iles				time zone	SYSTEM -
	Previous	Next	Finish K	Cancel	Compatibi <u>l</u> ity leve	I 5.2.003 ‡

FIGURE 6.2: Step 1 of import Configuration Wizard of the Read Excel operator.

68

The first step is choosing the Excel file that contains the required data. Step 2: Select the required sheet and cell range.

Data import wizard - Step 2 of 4 This wizard guides you to import your data. Step 2: An Excel file can contain multiple sheets. Please select the one you want to import into RapidMiner. Furthermore, you can mark a range of cells to be loaded.									
Nur	sery.xls								
Α	В	С	D	E	F	G	н	1	
usual	proper	complete		convenient	convenient		recommen	recommend	
usual	proper	complete		convenient	convenient		priority	priority	
usual	proper	complete		convenient	convenient		not_recom	not_recom	
usual	proper	complete		convenient	convenient	slightly_pr	recommen	recommend	
usual	proper	complete		convenient	convenient	slightly_pr	priority	priority	
usual	proper	complete		convenient	convenient	slightly_pr	not_recom	not_recom	
usual	proper	complete		convenient	convenient	problemat	recommen	priority	
usual	proper	complete		convenient	convenient	problemat	priority	priority	
usual	proper	complete		convenient	convenient	problemat	not_recom	not_recom	
	Previous Next Finish Cancel								

FIGURE 6.3: Step 2 of import Configuration Wizard of the Read Excel operator.

RapidMiner retrieves the Excel file and displays it in this step. The desired Excel sheet and cells are selected. By default, all cells are selected. In this process, there is no need to change anything at this step because there is only one sheet and all the required cells are already selected.

Step 3: Annotations.

This step enables the annotations to be specified. By default, the annotation of the first row is set to 'name'. If the first row of data does not contain the names of attributes, its annotation should be unchecked (as shown in Figure 6.4).

Step 4: Set the name, type, and role of attributes.

This is the most crucial step of the import wizard. The name, type, and role of attributes are specified in this step. Figure 6.5 shows the name, type, and role of attributes that are specified for the Nursery dataset. A breakpoint is inserted after this operator so that the output of the Read Excel operator can be seen in the Results Workspace.

6.3.2 Pre-processing

The Select attributes operator is used for selecting only the Form and Parents attributes. Only two attributes are selected to simplify this process (in order to make the explanation of working of Naïve Bayes operator simple).

6.3.3 Model Training, Testing, and Performance Evaluation

The Split Data operator is used in this process for splitting the dataset into two partitions. One partition is used for training the model and the other partition is used for testing the model. The *partitions* parameter of the Split Data operator performs almost the same task that the *split ratio* parameter does in the Split Validation operator. The *partitions*

70

Data import wizard - Step 3 of 4										
j.	This w Step 3 annota of the a proper assign	This wizard guides you to import your data. Step 3: In RapidMiner, each attribute can be annotated. The most important annotation of an attribute is its name - a row with this annotation defines the names of the attributes. If your data does not contain attribute names, do not set this property. If further annotations are contained in the rows of your data file, you can assign them here.								
Annotation	Α	в	С	D	E	F	G	н	1	
-	usual	proper	complete	1	convenient	convenient	nonprob	recomm	recomme	
-	usual	proper	complete	1	convenient	convenient	nonprob	priority	priority	
-	usual	proper	complete	1	convenient	convenient	nonprob	not_reco	not_recor	
-	usual	proper	complete	1	convenient	convenient	slightly_	recomm	recomme	
-	usual	proper	complete	1	convenient	convenient	slightly_	priority	priority	1
-	usual	proper	complete	1	convenient	convenient	slightly_	not_reco	not_recor	
-	usual	proper	complete	1	convenient	convenient	problem	recomm	priority	
-	usual	nroner	complete	1	convenient	convenient	nrohlem	nriority	nriority	
Previous Next Finish Cancel										

FIGURE 6.4: Step 3 of Import Configuration Wizard of the Read Excel operator.

Data impo	import wizard - Step 4 of 4 This wizard guides you to import your data. Step 4: RapidMiner uses strongly typed attributes. In this step, you can define the data types of your attributes. Furthermore, RapidMiner assigns roles to the attributes, defining what they can be used for by the individual operators. These roles can be also defined here. Finally, you can rename attributes or deselect them entirely.							
Reloa	Reload data							
\checkmark	\checkmark	\checkmark	\checkmark	\checkmark	\checkmark	\checkmark	\checkmark	\checkmark
Parents	Has_nur	Form	Children	Housing	Finance	Social	Health	Rank
polyno 💌	polyno 🔻	polyno 🔻	polyno 🔻	polyno 🔻	binomi 🔻	polyno 🔻	polyno 🔻	polyno 🔻
attribute 💌	attribute 💌	attribute 🔻	attribute 🔻	attribute 💌	attribute 💌	attribute 💌	attribute 🔻	label 🔻
usual	proper	complete	1	convenient	convenient	nonprob	recommend	recommend
<								Þ
📀 0 errors.						🖌 Ignore	errors 📃 S	Show only <u>e</u> rrors
Row	v, Column		Error		Original v	alue	Me	ssage
ext Einish Scancel								

FIGURE 6.5: Step 4 of import configuration Wizard of the Read Excel operator.

parameter defines the number of partitions and ratio of each partition. The *sampling type* parameter of the Split Data operator behaves exactly the same as the *sampling type* parameter in the Split Validation operator. The parameters of the Split Data operator are shown in Figure 6.6.

S Edit Parameter List: partitions	۶ 🔦	Split Data
Edit Parameter List: partitions The partitions that should be created.	partitions sampling type	Edit Enumeratio
ratio	use local random	seed
0.7		
0.3		
Add Entry		

FIGURE 6.6: The parameters of the Split Data operator.

Training the Model The Split Data operator provides the first partition (ratio = 0.7) through its first output port. This partition is used for training a Naïve Bayes classification model through the Naïve Bayes operator. Initially, all parameters of the Naïve Bayes operator are used with default values. The Naïve Bayes operator provides the Naïve Bayes classification model as its output. This model is given as input to the Apply Model operator.

Testing the Model The Split Data operator provides the second partition (ratio = 0.3) through its second output port. This partition is used for testing the Naïve Bayes classification model through the Apply Model operator. The Apply Model operator applies the Naïve Bayes classification model (that was provided at its *model* input port) on the testing dataset (that was provided at its *unlabeled data* port). The resultant labeled dataset is delivered as output by the Apply Model operator. This labeled dataset is provided as input to the Performance (Classification) operator.

Performance Evaluation The Performance (Classification) operator should be used for performance evaluation of classification models because it provides many different performance measures in its performance vector. The performance vector generated by the Performance (Classification) operator and the labeled data set are connected to the *result* port of the process.

6.3.4 A Deeper Look into the Naïve Bayes Algorithm

The results of this process show three outputs:

- 1. Performance Vector
- 2. Labeled dataset (obtained by application of model on testing dataset)
- 3. Naïve Bayes classification model

Table 6.1 shows the posterior probabilities of the label values. Posterior probability is calculated by dividing the number of examples with that label value by the total number of examples.

TIBBLE OUT I OBTORIOF PROBADINITION							
Rank (label)	No. of Examples	Posterior Probability					
recommend	1	1/9072					
priority	2986	2986/9072					
not_recom	3024	3024/9072					
very_recom	230	230/9072					
spec_prior	2831	2831/9072					

TABLE 6.1: Posterior probabilities

Figure 6.7 shows the first row of the labeled dataset. There is a confidence attribute for each possible value of the label. The label value with the highest confidence is assigned as predicted value for the example. In this example, the highest confidence (i.e., 0.450) is for label value = priority. Therefore, this example is predicted as priority.

The following steps explain how the confidences for each label value are calculated. They are calculated using the posterior probabilities and the distribution table (shown in Figure 6.8). The following calculations are done for finding the confidence of a label value:

- 1. Multiply the probabilities of all the attributes for that label value.
- 2. Multiply this product with the posterior probability of that label value.
- 3. Divide the resultant value by sum of all the confidences.

Here is an explanation of how the confidences for the first example were calculated. The rows where "Parents = usual" and "Form = complete" are highlighted in Figure 6.8 because these are the only rows that will be used for calculations for confidences in this example (because in this example, Parents and Form attributes have values "usual" and "complete" respectively). Firstly, confidences of all the label values are calculated without normalization. These values will not match the values in Figure 6.7. After calculating all the confidences, these values will be normalized. These normalized values will be the same as the values in Figure 6.7.

Rank	confidence(recommend)	confidence(priority)	confidence(not_recom)	confidence(very_recom)	confidence(spec_prior)	prediction(Rank)	Parents	Form
priority	0.001	0.450	0.339	0.062	0.148	priority	usual	complete

FIGURE 6.7: First row of labeled ExampleSet.

confidence (recommend) confidence (recommend) = P(Parents = usual |Rank = recommend) * P(Form = complete |Rank = recommend) * Posterior(recommend)

$$= 1.000 * 1.000 * 1/9072 = 0.0001$$

confidence (priority) confidence (priority) = P(Parents = usual |Rank = priority) * P(Form = complete |Rank = priority) * Posterior(priority)

= 0.449 * 0.262 * 2986/9072 = 0.0387

confidence (not_recom) confidence (not_recom) = P(Parents = usual |Rank = not_recom)
* P(Form = complete |Rank = not_recom) * Posterior(not_recom)

= 0.344 * 0.254 * 3024/9072 = 0.0291

O Text View	Text View Plot View Distribution Table Annotations						
Attribute	Parameter	recommend	priority	not_recom	very_recom	spec_prior	
Parents	value=usual	1.000	0.449	0.344	0.587	0.191	
Parents	value=pretentious	0.000	0.348	0.329	0.413	0.314	
Parents	value=great_pret	0.000	0.203	0.326	0.000	0.495	
Parents	value=unknown	0.000	0.000	0.000	0.000	0.000	
Form	value=complete	1.000	0.262	0.254	0.357	0.214	
Form	value=completed	0.000	0.256	0.244	0.300	0.245	
Form	value=incomplete	0.000	0.247	0.252	0.230	0.267	
Form	value=foster	0.000	0.235	0.250	0.113	0.274	
Form	value=unknown	0.000	0.000	0.000	0.000	0.000	

FIGURE 6.8: Distribution table.

confidence (very_recom) confidence (very_recom) = P(Parents = usual |Rank = very_recom) * P(Form = complete |Rank = very_recom) * Posterior(very_recom)

> 0.0053= 0.587230/9072 =0.357* *

confidence (spec_prior) confidence (spec_prior) = P(Parents = usual |Rank = spec_prior) * P(Form = complete |Rank = spec_prior) * Posterior(spec_prior)

> = 0.191 * 0.214 * 2831/9072 =0.0127

Normalization All the confidences are divided by the sum of all the confidence values to get normalized confidence values. The sum of confidence values is 0.0859. The normalized confidence values are shown in the table below.

TABLE 6.2 : N	ormalized confidences.
Rank (label)	Confidence
Recommend	0.0001/0.0859 = 0.001
Priority	0.0387/0.0859 = 0.450
not_recom	0.0291/0.0859 = 0.339
very_recom	0.0053/0.0859 = 0.062
spec_prior	0.0127/0.0859 = 0.148

As the confidence of the "priority" label value is the highest, the rank of this example is predicted as "priority".

Chapter 14

Robust Language Identification with RapidMiner: A Text Mining Use Case

Matko Bošnjak

University of Porto, Porto, Portugal; Rudjer Boskovic Institute, Zagreb, Croatia

Eduarda Mendes Rodrigues

University of Porto, Porto, Portugal

Luis Sarmento

Sapo.pt - Portugal Telecom, Lisbon, Portugal

Acrony	ms	213				
Introduction						
The Pr	oblem of Language Identification	214				
Text R	epresentation	217				
14.3.1	Encoding	217				
14.3.2	Token-based Representation	218				
14.3.3	Character-Based Representation	219				
14.3.4	Bag-of-Words Representation	219				
Classifi	cation Models	220				
Implen	nentation in RapidMiner	221				
14.5.1	Datasets	221				
14.5.2	Importing Data	223				
14.5.3	Frequent Words Model	225				
14.5.4	Character n-Grams Model	229				
14.5.5	Similarity-based Approach	232				
Applica	ation	234				
14.6.1	RapidAnalytics	234				
14.6.2	Web Page Language Identification	234				
Summa	ary	236				
Acknow	vledgment	237				
Glossar	ту	238				
Bibliog	raphy	238				
	Acrony Introdu The Pr Text R 14.3.1 14.3.2 14.3.3 14.3.4 Classifi Implem 14.5.1 14.5.2 14.5.3 14.5.4 14.5.5 Applica 14.6.1 14.6.2 Summa Acknow Glossan Bibliog	AcronymsIntroductionThe Problem of Language IdentificationText Representation14.3.1 Encoding14.3.2 Token-based Representation14.3.3 Character-Based Representation14.3.4 Bag-of-Words Representation14.3.5 Character-Based Representation14.3.6 Laracter-Based Representation14.3.7 DatasetsImplementation in RapidMiner14.5.1 Datasets14.5.2 Importing Data14.5.3 Frequent Words Model14.5.4 Character n-Grams Model14.5.5 Similarity-based ApproachApplication14.6.1 RapidAnalytics14.6.2 Web Page Language IdentificationSummaryAcknowledgmentGlossaryBibliography				

Acronyms

 ${\bf API}$ - Application Programming Interface

ETL - Extract, Transform and Load

HTTP - HyperText Transfer Protocol

k-NN - k Nearest Neighbours

NLP - Natural Language Processing

SVM - Support Vector Machines

TF-IDF - Term Frequency - Inverse Document Frequency

UTF-8 - Unicode Transformation Format – 8-bit

XML - eXtensible Markup Language

14.1 Introduction

Language identification, the process of determining the language of machine-readable text, is an important pre-processing step in many information retrieval and web mining tasks. For example, the application of natural language processing (NLP) methods may require prior language identification, if the language of the text at hand is unknown. In order to properly execute stemming, sentence tokenization or named entity recognition, we need to identify the language of the text to successfully apply appropriate language technologies. In addition, language identification is essential in machine translation tasks. Some text classification tasks, such as sentiment analysis in social media, may also require language identification for filtering content written in a specific language.

There are numerous proprietary and open-source solutions for language identification, including stand-alone solutions and APIs. However, proprietary solutions are usually costly, and APIs may cease to exist, while open-source solutions may require deeper understanding of the code for further modification and implementation. Moreover, most solutions are usually fixed on a pre-defined set of languages.

This chapter aims to provide a walk-through of language identification fundamentals by first introducing the theoretical background, followed by a step-by-step guide on model building, according to the standard practice of data mining. We will put emphasis on understanding the problem, and solving it in an effective and costless manner within the open source environment RapidMiner. We will become familiar with RapidMiner's "Text Mining Extension", and learn how to create several simple and fast workflows for language identification. The methods we will use are generic, and applicable to any dataset, with the appropriate pre-processing. In addition, we will learn how to use the implemented models in custom projects, by exporting them as web services using RapidAnalytics, and deploying them in custom applications. We will go through a case study of a web application for language identification of web pages, using RapidMiner's "Web Mining Extension". Special emphasis will be given to the optimization of the exported workflow, to enable faster execution, and harmonious integration in other applications.

The Problem of Language Identification 14.2

The task of language identification can be simply described as discerning the language of a given text segment. For example, given a set of six sentences in different languages:

How many languages do you speak well? Wie viele Sprachen können Sie gut sprechen? Combien de langues parles-tu bien? Πόσες γλώσσες μιλάς καλά; ¿Cuántas lenguas hablas bien? Koliko jezika govorite dobro? Quantas línguas você fala bem?¹

the goal is to identify the language of each sentence. This task may be somewhat trivial for humans. We might be able to easily discern the language and possibly identify some of them, even though we might not speak any of those languages. The first sentence is easy, since we are reading this book written in English. We may say that the second sentence "feels" German, or that the third one "sounds" French. The alphabet of the fourth sentence looks completely different from the others, so we might guess it "looks like" Greek, while the fifth sentence might "lead to" Spanish. We might find the sixth sentence somewhat tricky since it looks different from the previous ones. If we are vaguely familiar with the Slavic family of languages we might say it "relates to" that family, since the sentence is written in Croatian. Finally, although we might find the last sentence "similar to" Spanish, it is, in fact, written in Portuguese.

When identifying languages, we use our knowledge of languages, acquired either actively or passively. To design algorithms for achieving the same task, we first need to systematize the knowledge needed for language identification. There are several indicators we can rely on when identifying languages, without knowing those languages at all. These indicators are the following:

- Alphabet Differences in symbols a language is written in are the most obvious feature for language identification. Even a person unfamiliar with languages at hand can discern different alphabets, and can easily be taught to identify different scripts due to their specificity. This, of course, does not hold for languages written in the same alphabet. Examples of several different alphabets are given in Figure 14.1.
- **Typical words** Words can be a strong language indicator, whether those words are unique for the language, or the most frequent ones. Unique words, such as $fika^2$ in Swedish, or $saudade^3$ in Portuguese, are an excellent indicator of their corresponding languages. Nevertheless, they may occur infrequently in text. The most frequent words, like *niet*

215

¹Source: Tatoeba Project: Open, collaborative, multilingual dictionary of sentences: http://tatoeba.org, Sentence $n^{\circ}682402$.

²"Coffee-break"

³"A deep emotional state of nostalgic longing for an absent something or someone that one loves", source: Wikipedia

Latin	data mining
Greek	εξόρυξη δεδομένων
Cyrillic	добыча данных
Hebrew	כריית נתונים
Arabic	استخراج البيانات
Japanese	データマイニング
Chinese	数据挖掘

FIGURE 14.1: Examples of different alphabets.

in Dutch, ne in Slovene, or ikke in Norwegian⁴ are the next obvious choice, and we will discuss them further in the chapter.

- Accented letters Languages in Latin alphabets often contain various accented letters, which can be present in only several languages, or can even be unique to a language. These letters can be used for narrowing down the choice of languages in question. The downside of this approach is that accented letters may not at all be present in the text, depending on the encoding used or the style of communication. For example, in informal communication, people frequently omit accents. Examples of accented letters are \check{E} in Czech, \ddot{O} in German, Swedish, Icelandic, Turkish, etc., and \check{S} in Slavic and Baltic languages.
- **Special symbols** If a language uses special symbols, rarely present or not present at all in other languages, they can be used as a strong language indicator. A weakness of this indicator is that, similar to accented letters, these symbols can be used rarely, or not used at all. Examples of such symbols are the inverted question mark $\dot{\beta}$, and exclamation mark \dot{j} , used only in Spanish, the "Scharfes S", β used in German, \hbar used in Maltese, or semicolon ; used as a question mark in Greek, among others.
- Letter combinations Combinations of n consecutive letters found in substrings of words are called character n-grams, and are powerful language indicators. By finding these combinations we find digraphs, word beginnings and endings, prefixes, suffixes, or any other frequent letter combinations typical for a language. Examples of these combinations are the digraphs Lj, Nj, $D\tilde{z}$ in Croatian; Zs, Gy, Ny in Hungarian; typical word beginnings like *be-*, *re-*, *sh-* in English, and *des-*, *pr-*, *da-* in Portuguese; and typical word endings like *-ns, -ur, -ont* in French, and *-ung, -eit* in German.
- Word length Some languages tend to form long words by concatenating distinct words, suffixes, and prefixes into single words with a compound meaning. Such words can frequently be found in Finnish, German, and Hungarian, among others. An example of such a composite word is *Geschwindigkeitsbegrenzung*⁵ in German.
- **N-gram distribution** Distribution of character n-grams per language is one of the most powerful indicators used for language identification. Different languages exhibit different character n-gram distributions, starting with single letters, also called 1-grams or

⁴All of them mean "no".

⁵"Speed limit".

unigrams. An example of the distribution of unigrams is given in Figure 14.2. In this example we observe a much bigger usage of letter a in Finnish than in the other two languages, an almost non-existent usage of letter j in Italian, and far greater usage of letter z in Czech than in both Finnish and Italian.



FIGURE 14.2: Distributions of unigrams in Italian, Czech, and Finnish show clear differences in the use of distinct letters per language. In this example, accented letters are stripped of their accents, and non-alphabet symbols are ignored.

In the remainder of this chapter, we show how to implement several language identification methods using the most popular techniques and features, both in literature and in practice. Namely, we implement the frequent words, the character n-gram, and the similarity-based techniques. Although some researchers [1] argue that the problem of language identification has been resolved, issues are still ongoing [2, 3], from the language identification of web pages [4, 5], short queries [6, 7], and micro-blog messages [8], to differentiation of similar languages [8, 9]. However, the extent of these issues go beyond the scope and the intention of this chapter. We continue by presenting the fundamentals of text representation and classification models.

14.3 Text Representation

In order to start the text analysis process, we first need to represent the text as a vector of numbers in a consistent way. Prior to explaining the representation format, we need a word of notice regarding digital alphabet representations — encodings.

14.3.1 Encoding

Simply put, character encoding is a convention that maps a number to a symbol, which enables numerical storage of characters on a digital media. Due to historical reasons, there are many different ways to represent different characters. There are two main ways of encoding characters: single- and multi-byte character encodings.

Single-byte character encodings represent symbols in a single-byte of information. The most known such encoding is ASCII, which uses only 7 bits to encode up to 128 different symbols. The inclusion of the 8th bit already brings many troubles since there are many different character encoding tables which will map numbers higher than 127 to different symbols. The most widespread single-byte character encodings are the ISO 8859, and Microsoft Windows sets of encodings. Due to single-byte encodings' limited capacity of only 256 symbols, multi-byte encodings are more frequent in use today. Multi-byte character en-

codings utilize several bytes of data to represent a symbol, and are therefore able to encode much larger spectra of symbols. The most widespread multi-byte encoding is UTF-8. 6

The choice of encoding is very important since a badly chosen encoding might make the text illegible. For example, the Greek word for data, $\delta \varepsilon \delta \delta \mu \varepsilon \nu \alpha$ encoded in UTF-8, if mistakenly shown with the ISO 8859-1 encoding, will appear as $\hat{I}'' \hat{I} \mu \hat{I} \hat{I}^{2} \hat{I}^{1}/4 \hat{I} \hat{I}^{1}/2 \hat{I} \pm$.

Throughout the rest of this chapter, we use the UTF-8 encoding, due to its widespread usage and the availability of major language scripts in it. While some language identification methods operate on bytes of data rather than characters, they cannot cope with code overlap in different encodings, and need to employ encoding detection. More information about encoding detection for language identification can be found in [10].

In addition, we limit our language identification to languages written in the Latin alphabet. Other alphabets can be detected by checking the availability of letters in certain code ranges, under the assumption of single encoding like the UTF-8. For example, we can detect the Greek alphabet by executing the regular expression $[\u0370-\u03FF]$ *. This regular expression will include all the symbols ranging from the code point 0370 to the code point 03FF, expressed in hexadecimal, which in UTF-8 represents all the new and old Greek alphabet symbols. In RapidMiner, this regular expression can be applied to text with the *Keep Document Parts* operator. Note that this regular expression will not include Latin letters nor numerals. To include them and the standard punctuation, $[\u0370-\u03FF]$ would be used.

With the encoding explained, and set to UTF-8, we continue on the level of text representation for data mining tasks.

14.3.2 Token-based Representation

The token-based representation is built on top of basic meaningful elements of text, called tokens. Tokens can be words separated by delimiters; logograms—signs and characters representing words or phrases, such as Chinese letters; idioms, or fixed expressions, such as named entities—names of places, people, companies, organisms, etc.

The process of token extraction is called tokenization. For most of the languages and uses, tokenization is done by splitting sentences over whitespaces and punctuation characters, sometimes ignoring numbers, or allowing specific punctuation in a word, e.g., *wasn't*, *off-line*. For example, the following sentence:

El mundo de hoy no tiene sentido, así que ¿por qué debería pintar cuadros que lo tuvieran? 7

separated by non-alphabet characters, results in the following set of tokens:

El mundo de hoy no tiene sentido así que por qué debería pintar cuadros que lo tuvieran.

In this chapter we use this simple tokenization over non-letter characters. However, in languages in which words are not delimited by space, such as Chinese, Japanese, Thai, Khmer, etc., tokenization is a more complicated process, and it requires complex machine learning models outside of the scope of this chapter.

Note that when tokenizing over non-letter characters, apostrophes are used as points of separation, which results in wasn't tokenized as wasn and t. We can argue that this is not

⁶Usage Statistics of Character Encodings for Websites http://w3techs.com/technologies/overview/character_encoding/all/.

⁷". Today's world doesn't make any sense, so why should I paint pictures that do?" --- Pablo Picasso

a desired behaviour in our process, since *wasn't* should be one token. However, due to the simplicity of setting parameters in RapidMiner, we opted for this solution. Easy tweaking of RapidMiner options allows inclusion of any other symbol, including the apostrophe. Later in the chapter, we explain how to set these parameters.

When dealing with specialized forms of text, like blog or micro-blog messages, for specific purposes like sentiment analysis, tokenization becomes a more complicated process in which non-letter tokens, such as smileys and interjections,⁸ are retained. For more information on tokenization, consult further literature [11, 12, 13].

Other than splitting the text in tokens, and continuing to build a data structure for text mining, we can also split the text on the character basis. This results in a character-based representation.

14.3.3 Character-Based Representation

Character-based representation is a text representation based on top-of-character ngrams, previously defined as word substrings consisting of n consecutive characters. In general NLP usage, n-grams can also denote an ordered set of words, however is this chapter when mentioning the term n-gram, we specifically refer to character n-grams. For example, the following sentence:

Ce n'est pas assez d'avoir l'esprit bon, mais le principal est de l'appliquer bien. 9

cut into 2-grams, also called bigrams, after tokenization results in the following set:

_a(3) _b(2) _c _d(2) _e(3) _l(3) _m _n _p(2) ai al ap as(2) av bi bo ce ci d_ de e_(3) en er es(3) ez ie in ip iq ir is it l_(3) le li ma n_(3) nc oi on pa(2) pl pp pr(2) qu r_(2) ri(2) s_(2) se sp ss st(2) t_(3) ue vo z_¹⁰

Extraction of n-grams can be executed after tokenized, or even non-tokenized text. The difference is that when extracting n-grams from non-tokenized texts, depending on the number n, n-grams can catch both the ending of one, and the beginning of another word. For example, when extracting 3-grams from the term *data mining*, 3-gram "a m" will also be extracted. This does not happen if we execute tokenization prior to n-gram extraction.

N-grams are normally used to catch word beginnings and endings, as well as typical combinations of consecutive letters. Whereas the number of words is practically unlimited, the number of n-grams is limited to the power of number of letters n. In real scenarios, the number of observed n-grams is actually much lower than the maximum possible number of n-grams, since not all letter combinations are present in natural languages. This enables compressed format of the text and creates helpful features in language identification.

Whether we choose to represent the text with tokens or n-grams, prior to continuing the text mining process, we convert them to a bag-of-words vector.

14.3.4 Bag-of-Words Representation

The bag-of-words is a representation used in NLP and text mining, in which a text, such as a sentence or a full document, is represented by a set of words—an unordered collection of

 $^{^{8}}$ Words used to express emotions and sentiments.

⁹"It's not enough to have a good mind; the most important thing is to use it well." ---René Descartes. ¹⁰Apostrophes replaced by empty space, added underscores to empty spaces to emphasize word begin-

nings and ends shorter than n, bigram cardinality shown in parenthesis, where greater than one.

words.¹¹ Note that although this representation completely disregards grammar and word order, it suffices for our needs. For example, the following sentences:

The following sentence is true. The preceding sentence is false.

consists of the following words, without repetition:

1. the, 2. following, 3. sentence, 4. is, 5. true, 6. preceding, 7. false

By using this set of words as an index or dictionary, we create a fixed representation of a sentence, consisting of an occurrence of a specific word, at its appropriate index. This results in creating a word vector. For example, the previous two sentences, using the extracted dictionary, result in these word vectors:

 $(1 \ 1 \ 1 \ 1 \ 1 \ 0 \ 0)$ $(1 \ 0 \ 1 \ 1 \ 0 \ 1 \ 1)^{12}$

The second vector in our example denotes a sentence containing one element of each of indexes 1, 3, 4, 6, 7 in the dictionary, corresponding to words *the*, *sentence*, *is*, *preceding*, and *false*. With this approach we lose grammatical structures and word order, but we simplify the resulting model, making it tractable from the analytic point of view, and as we will demonstrate, sufficient for our case studies. Upon obtaining word vectors of the bag-of-words representation, we are ready to use classifiers for language identification.

14.4 Classification Models

Our main goal, developing a system for language identification, requires an applicable classification model. An example of a model might be a simple string matching—if a language contains one or more words, or n-grams of a certain language, it belongs to that language. This type of model leads to rule-based systems. Another approach is to represent the text segment in an *n*-dimensional space, where each dimension corresponds to a single word metric like the occurrence, TF-IDF score, etc. This leads to vector-space models, that are frequently used with Support Vector Machines (SVMs). In practice, statistical language models, like Naïve Bayes or Markov models, are frequently used since they yield high performance, and in the case of Naïve Bayes, low computational cost. Though the Naïve Bayes classifier is often outperformed by other classifiers, like SVM, in terms of prediction performance, we use it since it is fast, and does not need parameter optimization. Feel free to experiment with other classifiers to find the best performing one.

The first classification model we employ is the Naïve Bayes classifier. The Naïve Bayes classifier is a generative probabilistic model that assigns a probability to a segment of a text, composed of either single or multiple tokens or parts of tokens [14]. This is done through statistical analysis of the text, by translating frequencies of word occurrences in a document into probabilities, and employing a statistical framework for classification [15]. This is one of the models we use in our case studies.

The second model we employ is the similarity-based model. The similarity-based model

¹¹The precise term would be multiset, a generalization of set, which allows repeating members.

¹²The example is considered case-insensitive; words are ordered in the sentence order.

Robust Language Identification with RapidMiner: A Text Mining Use Case 221

relies on creating a distribution of elements, such as tokens or n-grams per language, and using a similarity measure to determine the language with the most similar distribution to the queried example. While the classification approach requires a large number of examples per language, the similarity approach relies on a single large example per language, for example, all the texts of a single language concatenated together. This model, in its essence, is the k Nearest Neighbours (k-NN) model. Having a large text per language enables us to extract a distribution of tokens, or n-grams, that closely approximates the true distribution of the language.

One of the most widely used approaches is the ranking approach proposed by Cavnar and Trenkle [16] in which they implement the ranking similarity between the query text profile and language profiles. Since the ranking similarity is not yet implemented in RapidMiner, we can use other similarity measures frequently utilized, such as the Euclidean distance, dot product, and cosine similarity [10, 17, 18].

Having token and character-based text representation as the ground representation, out-of-bag word vectors, and the Naïve Bayes and k-NN classifiers at hand, we are in the position to start discussing the implementation of possible language identification workflows in RapidMiner.

14.5 Implementation in RapidMiner

In this section, we move on to building RapidMiner workflows by presenting three different models for language identification, built on top of the previously introduced theoretical background. The models we present are the frequent words, character n-gram, and the similarity-based models [19]. Of course, additional operators in RapidMiner can be used to build more complex models, and improve the ones we will present; however, we do not cover all the advanced details in this chapter.

The first step in the implementation process is to obtain a suitable dataset for experimenting with language identification techniques. If you already have a dataset available, and you wish to skip on to model construction, you can proceed to Section 14.5.3. However, if you do not have a dataset prepared, and cannot easily get a hold of one, there are several freely available datasets you can use. In the next subsection we introduce a few such datasets, which you can use to re-create and experiment with RapidMiner workflows.

14.5.1 Datasets

When creating a dataset for developing language identification models, one has to consider several important factors:

- **Time** Language is dynamic, and changes over time. Commonly used words and phrases differ between centuries, decades, and even years. Obtaining datasets from comparable eras ensures similar language properties and stability of used words.
- **Source** Different sources of text exhibit different properties. For example, news articles, web pages, dictionaries, social media like Twitter, and classical writing greatly differ in the vocabulary, punctuation, and even word case. Thus, the choice of dataset should favor texts from languages coming from the same or comparable sources. This ensures that the language has a fairly stable style.

Size Datasets may differ in size. Small datasets might not exhibit a very representative

distribution of words or n-grams, while a large dataset might. However, computations with large datasets are more costly, in terms of time, complexity, and memory requirements. Obtaining comparable sizes of datasets per language is also important to represent all the languages with approximately the same number of sentences, words, or tokens.

Length Length of texts used for training or for querying has an impact on performance. Building classifiers on small sentences for training might decrease performance if the number of sentences is also small. Likewise, identifying language of a short sentence is a difficult task since only a small number of features are observed.

Therefore, when creating a dataset, we should collect texts of comparable sizes, coming from comparable sources, created in near time, and of adequate text length compared to query texts, in order to ensure higher performance of the language identifier.

There are multiple different datasets which can be used for this purpose, that abide by most of the previous recommendations.

- Wikipedia Wikipedia¹³ is a free Internet encyclopaedia with more than 21 million articles written in 284 languages. The articles can be downloaded and used for language identification. Wikipedia database dumps can be found at http://dumps.wikimedia.org/. in various formats and levels of information.¹⁴ For language identification tasks, we recommend abstracts of web pages in XML format.
- Leipzig Corpora Collection The Leipzig Corpora Collection¹⁵ is a collection of textual corpora in 155 different languages using the same format and comparable sources [20]. The corpora consist of randomly selected sentences in the corpus language, and are available in sizes varying from tens of thousands up to tens of millions of sentences. The sources of sentences are news articles and texts, randomly collected from the web.
- **Project Gutenberg** Project Gutenberg¹⁶ is the oldest digital library of written cultural work. This web page offers full text of public domain books in various formats, out of which UTF-8 plain-text is the most straightforward for our use. In May 2012, Project Gutenberg claimed to have in excess of 39,000 free ebooks in more than 50 languages, with English being the most predominant language. Note that these books are older belletristic works.
- **European legislature datasets** There are several datasets originating from the official documents of the European Union (EU), published in the official languages of the EU. Download

 $EuroGOV^{17}$ is the document collection of web documents crawled from European governmental sites, with restricted access.

European Parliament Proceedings Parallel Corpus¹⁸ [21] is a parallel corpus extracted from the proceedings of the European Parliament 1996–2010. It consists of 21 languages and 20 language pairs, English versus the rest. Though conceived as a standard dataset for statistical machine translation systems, it can also be used for language identification.

¹³http://www.wikipedia.org/.

¹⁴Details on database download are available at http://en.wikipedia.org/wiki/Wikipedia:Database_download.

¹⁵http://corpora.informatik.uni-leipzig.de/download.html.

¹⁶http://www.gutenberg.org/.

¹⁷http://ilps.science.uva.nl/WebCLEF/EuroGOV/.

¹⁸http://www.statmt.org/europarl/.

The JRC-Acquis Multilingual Parallel Corpus¹⁹ is the corpus representing the total body of European Union law applicable in the EU member states [22]. These texts are available in 22 languages.

In this chapter, we use the Leipzig Corpora Collection dataset of 30.000 sentences in English, German, French, Spanish, and Portuguese. The dataset is easy to download, we just need go to the Leipzig Corpora Collection web page at http://corpora.informatik.uni-leipzig.de/download.html, scroll down to "Download Corpora" type in the captcha text, and click on "check". When a table with various corpus sizes and formats appears, we select the corpora we find suitable and download them by clicking on them. The next step after downloading the data is to prepare it and import into RapidMiner.

14.5.2 Importing Data

Upon downloading the corpora for English, German, French, Spanish, and Portuguese, we extract the compressed file contents and isolate files with a suffix -sentences.txt. These sentence files are tab delimited files containing a line number and a sentence in each line of the file. After downloading a dataset per language, and extracting its content, we rename the downloaded text file according to the language of the data, in our case *eng*, *deu*, *fra*, *spa* and *por* respectively, and save the files in a single directory. In our repository, these files are located in the data/corpus directory.

In order to import the Leipzig Corpora Collection, we create a data importing workflow, presented in the upper part of Figure 14.3. We loop over the files we previously saved using the *Loop Files* operator (1), by setting the *directory* property to the path to data, in our case data/corpus. The *Loop Files* operator will iterate over all files in the directory we set it to, and execute its nested process for each iterated file.

The content of the *Loop Files* operator is shown in the lower part of Figure 14.3. Essentially, in every iteration of the *Loop Files* operator, we read the text file, and fit it to specifications for further processing, creating the appropriate label, and finally splitting the data into two datasets: training and test datasets. The training dataset is further split into two forms: i) example per sentence and ii) all the text concatenated per language.

In more detail, the *Read* CSV operator (1) receives the file to open from the nested process of the Loop Files operator. We set the column separator property to t since the text is tab delimited, uncheck the use quotes property since sentences are not delimited by quotes, uncheck the *first row as names* property since there are no row names in the first row of the file, and set the *encoding* property to UTF-8. The *Rename* operator (2) renames the default att2 attribute to text, since the *Read CSV* operator returns default values for row names, when the first row of the file does not convey names of attributes. With the Select Attributes operator (3), we select only the text attribute, thus ignoring the line number. The *Generate Attributes* operator (4) generates the language attribute, and sets its value to the replace("%{file_name}",".txt","") expression. This expression essentially uses the name of the file to populate the language attribute. This is why we renamed the file names to short names of languages, to automate the whole process. The replace("%{file_name}",".txt","") expression relies on the macro %{file_name} constructed by the *Loop Files* operator per each iteration, and contains the name of the file in the current iteration. The next step is to set the label role to the newly created 'language' attribute with the Set Role operator (5). The output of the Set Role operator is now a dataset containing two attributes: the *text* attribute containing each sentence text, and the language attribute containing the label of each of those sentences, fixed to the name of

223

¹⁹http://langtech.jrc.it/JRC-Acquis.html.



FIGURE 14.3: The workflow for loading the Leipzig Corpora dataset. The upper image depicts the main workflow process, whereas the lower image depicts the content of the *Loop Files* operator (1) in the upper image.

the file which is a short code for the language. With the *Split Data* operator (6) we split the dataset randomly with the ratios of 0.3 and 0.7 into the train and the test datasets, respectively. This kind of split is unusual in real-world setting—we would usually create a bigger training set than the test set. However in this case, we opted for a smaller training set to increase the speed of training. We output the test dataset without any change, while we process the train set into two modes. The first mode is simply the standard sentenceper-example dataset, outputted for processing and saving, just as the test set. The second mode is the concatenated dataset, created by first converting the text attribute from nominal to text type, using the *Nominal to Text* operator (8), since the nominal type cannot be processed as a textual value. In the end, the *Data to Documents* (9) and the *Combine Documents* (10) operators convert the dataset first to a collection of documents, and then combine that collection in a single resulting document.

Back to the main process. The datasets are outputted from the *Loop Files* operator (1) in the following order: the collection of documents containing one document of concatenated train data, the collection of train data, and the collection of test data, all three per iterated file. The collection of documents containing documents of concatenated data per language is converted to data with the *Documents to Data* operator (2), the role of the *language* attribute is set to label via the Set Role operator (3). The resulting dataset is the concatenated dataset, a dataset containing two columns: a label column denoting the language, and the text column, containing concatenated texts of all the sentences corresponding to each language. This dataset is stored to the repository as the language_concatenated dataset, using the *Store* operator (4). The collection of train and test data, on the other hand, both go through the identical process. They are first flattened in a single dataset by the Append operator (5), and the text attribute is converted to the text type via the Nominal to Text operator (6). The resulting datasets contain two columns: a label column denoting the language, and the text column, containing one sentence per line of the dataset. These datasets are stored in the repository as language_train and language_test with the Store operators (7) and (10), respectively.

Thus far, we obtained the data for language modeling, and we are ready to construct the first language identification system. The simplest and the most naïve approach to language identification is the dictionary-based identification. By enumerating all the possible words in every single language, we can try to match the word with its language and identify the language of the text. However, this method is intractable, since the language is a dynamic complex system with no bound on the number of words. In addition to the complexity, the creation and maintenance of a large dictionary set would create another set of problems, from computational load, to the fact that different languages share same words. Alternatively, we can select a fixed number of words and use them for language identification. Choosing the most frequent words for each language is a reasonable option. This is why the next model we introduce is the frequent words model.

14.5.3 Frequent Words Model

The frequent words approach identifies the most frequent words in the text, and uses them for the language identification. The most used words across different languages have similar meanings, and regarding the type of words, they are usually articles, prepositions, pronouns, and some frequently used verbs, and adjectives. Due to their high frequency of occurrence, they are suitable even for shorter texts. For example, it might be easy to see that the following sentences are written in three distinct languages:

225

Det gäller att smida medan järnet är varmt. 20 Het is niet alles goud wat er blinkt. 21 Nu masura pe altii cu palma ta. 22

Nevertheless, it is difficult to identify specifically the languages of these sentences without being acquainted with them, especially if those languages are not widespread. Still, given a list of the most frequent words in Dutch,²³ Romanian, ²⁴ and, Swedish²⁵ as in Table 14.1, it is straightforward to identify the language of each of these sentences as Swedish, Dutch, and Romanian, respectively.

прпп	The house of the 20 most nequent words in Daton, reemainan, and Sweda
Dutch	ik, je, het, de, dat, is, een, niet, en, wat, van, we, in, ze, hij, op, te, zijn, er,
	maar
Romanian	sa, nu, de, o, ca, și, a, ce, în, e, am, pe, te, la, mai, cu, un, asta, ai, ma
Swedish	i, och, en, av, som, är, att, den, på, till, med, för, det, han, de, ett, var, har,
	från, under

TABLE 14.1: Lists of the 20 most frequent words in Dutch, Romanian, and Swedish.

Building a statistical model for language identification using the frequent words method in RapidMiner is straightforward. First, we must analyze the training set text, extract all the words from it, and select the most frequent ones. Second, we represent the analyzed text by the most frequent words, converting the text to the word-vector representation in which a value of an attribute denotes the occurrence or a frequency measure of a given word in the text. Finally, we feed the word-vector representation of the text to the Naïve Bayes classifier, and build a model that can classify languages by the most frequent words exhibited in the training set.

We create a single RapidMiner workflow for the purpose of building the frequent words language identification, evaluating it, and applying it on a test dataset. This workflow is presented in Figure 14.4.

The first three operators (1), (2), and (3) execute data loading, textual processing, and model training and performance estimation via 10-fold cross-validation. First we load the language_train dataset using the *Retrieve* operator (1). The dataset we load contains two attributes, the *text* attribute containing sentences, and the *language* attribute containing the language of the text as the label. After loading the data, we analyze it with the *Process Documents from Data* operator (2). This operator created word vectors based on term frequencies with the property *Vector creation* set to Term Frequency. The inner Vector Creation Subprocess of the Process Documents from Data operator, as can be observed in the middle workflow of the Figure 14.4, contains the *Transform Cases* and the *Tokenize* operators. These two operators are set to the default properties; therefore each sentence of our dataset is first transformed to lower-case, and then tokenized over non-letter characters. We want to limit the resulting number of words in the analysis, to speed up the model learning, as well as simplifying the process. We choose to extract the top 0.1% of all the words in our example set by setting the *prune method* parameter to ranking, with the

²⁰"You should forge while the iron is hot."

²¹"All that glitters is not gold."

²²"Don't measure others with your hand."

 $^{^{23} \}tt http://en.wiktionary.org/wiki/Wiktionary:Frequency_lists/Dutch_wordlist.$

²⁴http://en.wiktionary.org/wiki/Wiktionary:Frequency_lists/Romanian.

 $^{^{25} \}tt http://www.lexiteria.com/word_frequency/swedish_word_frequency_list.html$



FIGURE 14.4: Workflow for the language identification system based on the most frequent words in languages. The main process is depicted in the top image. The middle image depicts the contents of both *Process Documents from Data* operators (2), (5). The bottom image depicts the content of the *Validation* operator (3).

prune below ranking parameter set to 0.001, and the prune above ranking set to 0.0. The final stage of the modeling phase is the X-Validation operator (3), which executes the Naïve Bayes classifier, and estimates its performance through 10 folds.

The next four operators (4), (5), (6), and (7), constitute the application phase of the workflow, meaning that they are responsible for loading the test data, transforming it in the way in which the previously built model was trained, applying the model, and calculating its performance. First, we load the language_test dataset using the *Retrieve* operator (4). The test dataset is of the same format as the train dataset. After loading the train dataset, we analyze it in the *Process Documents from Data* operator (5) by passing the WordList from the *Process Documents from Data* operator (2). This ensures that the same words are used for analyzing the test set, and thus ensuring that the dataset is of the same format as the dataset used for model building. The operators inside the *Process Documents from Data* operator (5) are the same as in the model building phase, namely the *Transform Cases* and the *Tokenize* operators. The result of the *Process Documents from Data* operator is a word-vector dataset. The *Apply model* operator (6) then receives the Naïve Bayes model, outputted by the *Validation* operator (3), and applies it on the newly created word-vector dataset. The process ends with the display of the performance.

Note that if we want to separate the model building and the model application into distinct workflows, we should save the resulting WordList, and the model, both built during the model building phase, and load them in the application phase of the workflow of interest.

The cross-validation estimation of the accuracy of the frequent words workflow is equal to $96.19\% \pm 0.35\%$, with the 96.23% accuracy achieved on the test set, as shown in Figure 14.5. It is possible to further improve these results by increasing the percentage of the words included in the model, however this would increase the model complexity. For example, including the top 1% of words instead of the top 0.1%, results in the 99.05% accuracy on the test set.

Table View	O Plot View									
accuracy: 96.19% +/- 0.35% (mikro: 96.19%)										
	true deu	true eng	true fra	true por	true spa	class precision				
pred. deu	8575	11	16	6	0	99.62%				
pred. eng	62	8736	8	3	10	99.06%				
pred. fra	13	22	8749	5	35	99.15%				
pred. por	10	17	18	8324	53	98.84%				
pred. spa	340	214	209	662	8902	86.20%				
class recall	95.28%	97.07%	97.21%	92.49%	98.91%					
Table View	O Plot View									
Table View accuracy: 96.2	Plot View									
Table View accuracy: 96.2	Plot View 23% true deu	true eng	true fra	true por	true spa	class precision				
Table View accuracy: 96.2 pred. deu	Plot View 23% true deu 20013	true eng 28	true fra 27	true por 18	true spa 3	class precision 99.62%				
Table View accuracy: 96.2 pred. deu pred. eng	Plot View Plot View true deu 20013 152	true eng 28 20436	true fra 27 11	true por 18 12	true spa 3 22	class precision 99.62% 99.05%				
Table View accuracy: 96.2 pred. deu pred. eng pred. fra	Plot View Plot View true deu 20013 152 25	true eng 28 20436 38	true fra 27 11 20402	true por 18 12 30	true spa 3 22 77	class precision 99.62% 99.05% 99.17%				
Table View accuracy: 96.2 pred. deu pred. eng pred. fra pred. por	Plot View 23% true deu 20013 152 25 35	true eng 28 20436 38 50	true fra 27 11 20402 39	true por 18 12 30 19403	true spa 3 22 77 113	class precision 99.62% 99.05% 99.17% 98.79%				
Table View accuracy: 96.2 pred. deu pred. eng pred. fra pred. por pred. spa	Plot View 23% true deu 20013 152 25 35 775	true eng 28 20436 38 50 448	true fra 27 11 20402 39 521	true por 18 12 30 19403 1537	true spa 3 22 77 113 20785	class precision 99.62% 99.05% 99.17% 98.79% 86.37%				

FIGURE 14.5: Estimated performance of the frequent words method with 10-fold cross-validation (upper image) and performance of the method on the test set (lower image).

Taking a look at the top 10 words extracted from the training set gives us a glimpse of the dataset, and can provide a quick check if everything is in order. This is done by taking a look at the WordList output of the *Process Documents from Data* operator (2), and sorting it over *Total Occurrences*, as presented in Figure 14.6. Most of the top words are, by themselves, bad identifiers of a single language, since they are shared among different

Robust Language Identification with RapidMiner: A Text Mining Use Case

229

languages. In our example, *a* is shared among all languages analyzed, except German, and *en* is shared among French and Spanish, though *the* is specific only to the English language. We also note that most of the top 10 words do appear in all languages, since it is possible for a sentence in, for example, French to contain foreign named entities such as "El País" or "The New York Times". Though most of the top words are not reliable indicators of a language by themselves, when taken together they achieve a higher accuracy. Note here that the most frequent words are usually discarded in various text mining applications on a single language since those words appear often, and as such do not have the needed discriminating value for various other purposes, though they do in the case of language identification.

Word	Attribute Name	Total Occurences 🔻	Document Occurences	deu	eng	fra	por	spa
de	de	27488	15832	30	17	10193	6809	10439
а	а	13553	10592	31	4116	1066	4935	3405
la	la	12868	8542	4	11	5891	39	6923
the	the	9629	5606	13	9525	77	6	8
que	que	9565	7429	0	0	1048	4066	4451
en	en	7372	5738	2	6	2897	1	4466
in	in	5502	4644	2335	3138	18	6	5
el	el	5462	3892	5	6	19	6	5426
0	0	4995	3899	5	19	23	4672	276
le	le	4884	3706	5	1	4605	3	270

FIGURE 14.6: List of the top 10 most frequent words in the language identification dataset, obtained with the frequent words workflow.

The advantages of the frequent words method are its speed, due to relatively small set of words, and high performance, greater than 95%. Higher accuracy can be achieved by including more words, thus increasing the model complexity and its processing time, or by using the two methods we present next. The disadvantages of the frequent words method are problems with similar languages due to similar top words, and bad results on sentences which do not include the top frequency words—usually short words and phrases.

A variation of this technique, short frequent words, uses only words of short length, usually four or five letters long. The rationale of the short frequent words technique is the same, with added shortness criteria which additionally filters out longer words, thus limiting itself to a smaller subset of the most frequent words. This filters out longer words which might be frequent due to linguistic properties of the source of the dataset. Creating a workflow for this technique is simple; it is only necessary to add a *Filter Tokens (by Length)* operator in the *Process Documents from Data* operator in Figure 14.4 (upper image), just after the *Tokenization* operator. The usual length of the token, set for this method, is 5 to 6 characters.

In the following subsection, we present the next method, which differs from the frequent words model on the level of the unit of text used for modeling; it uses n-grams rather than words. This is the character n-gram method.

14.5.4 Character n-Grams Model

We observed that the frequent words method achieves high accuracy, although it performs poorly in cases when the text to be classified does not contain any of the most frequent words. In cases of short sentences, or single words, we want a method capable of achieving high accuracy, without using words as the basis for classification. This is why we turn to character n-grams. As mentioned before, character n-grams are sequences of n consecutive characters, and as such are a powerful feature for language identification. They dissect the

text in its substrings and can catch beginnings and endings of words, as well as some typical character combinations, characteristic to a language. For example, the following sentences:

I cannot make out the meaning of this sentence. Ich kann die Bedeutung dieses Satzes nicht erkennen. 26

are easily classified into English and German, respectively, knowing that the character ngrams I, can, the, -ing, and of appear more often in English, and ich, die, be-, -ung, -cht in German. Notice that n-grams catch not only typical word beginnings and endings, or the combinations of letters, but also the short frequent words of length less than or equal to n.

Building a character n-gram language identification model is essentially an upgrade of the frequent words model, with the resulting workflow differing only slightly from the workflow in Figure 14.4. The conceptual difference is only in employing character n-gram creation after the tokenization, when analyzing the text. The rest of the process is the same: we convert the text to the word-vector representation, and build a classifier to classify it.

The workflow of the character n-gram method, depicted in Figure 14.7, is divided into two parts, i) the model creation consisting of operators (1), (2), and (3) and ii) the model application consisting of operators (4), (5), (6), and (7). As stated, the difference between the character n-gram and the frequent words workflows is minor. In the model building phase, the difference is in the two new operators in the *Process Documents from Data* operator (2), and a modification in properties of two other operators in it. The *Replace Tokens* operator, as seen in the lower image of Figure 14.7, adds underscores around a word, which enables differing cases when an n-gram is found at the beginning or at the end of a word. This is done by replacing the expression $([\w]+)$ with the expression ... This expression will replace any string of word characters longer than one, with the same string surrounded by three underscores from both the left and the right. The Generate n-Grams (Characters) operator creates character n-grams, in our case, n-grams of length 4. Note here that should we tick the *keep terms* option, we would keep the words, in addition to character n-grams, and thus effectively create a hybrid between the frequent words and the character n-gram methods. Two important modifications are apparent in the Process Documents from Data operator (2). First, the *prune method* property is set to by ranking, pruning all n-grams which occur in more than 10%, and less than 1% of sentences. Second, the *mode* of the tokenization is set to regular expression with the *expression* equal to [^a-zA-Z_]+. This regular expression describes all the splitting points. In this particular case it defines that the splitting points are everything NOT ([^]) containing strings of all the lower- and upper-case letters, and the underscore $(a-zA-Z_{-})$, of string length equal to at least one character ([]+). By adding additional characters or character ranges in this expression, we influence tokens used for the tokenization. For example, by adding 0-9 to it, we enable numbers and words with numbers and tokens, by adding ' we include apostrophes in words, and by adding various other punctuation characters we can catch smileys and interjections. The rest of the workflow stays the same as the workflow for the most frequent words model.

When executed, the character n-gram workflow, depicted in Figure 14.7, achieves $98.26\% \pm 0.14\%$ accuracy of the estimated performance with 98.31% accuracy on the test set, as observed in Figure 14.8. By increasing the number of n-grams included in the analysis, we improve the accuracy of the method.

When compared to the frequent words method, character n-grams exhibit an important property—they create a compact representation of text into a finite number of n-grams. This representation is a priori limited by the number of n-grams, whereas the word count in a language is not limited at all, due to the fact that new words are being produced constantly, i.e., language is dynamic.

²⁶Both of these sentences have the same meaning. Source: Tatoeba Project.


FIGURE 14.7: The workflow for the language identification system based on character ngrams depicted in the upper image. The lower image depicts the content of both *Process Documents from Data* operators (2), (5). The content of the *Validation* operator (3) is the same as in Figure 14.4.

Table View	O Plot View					
accuracy: 98.2	26% +/- 0.14% (mi	kro: 98.26%)				
	true deu	true eng	true fra	true por	true spa	class precision
pred. deu	8965	14	11	7	32	99.29%
pred. eng	5	8906	18	10	32	99.28%
pred. fra	24	64	8929	32	87	97.73%
pred. por	4	6	15	8764	196	97.54%
pred. spa	2	10	27	187	8653	97.45%
class recall	99.61%	98.96%	99.21%	97.38%	96.14%	
Table View	O Plot View					
accuracy: 98.3	31%					
	true deu	true eng	true fra	true por	true spa	class precision
pred. deu	20923	20	28	14	75	99.35%
pred. eng	18	20823	56	17	88	99.15%
pred. fra	42	118	20832	62	205	97.99%
pred. por	7	19	19	20475	458	97.60%
pred. spa	10	20	65	432	20174	97.45%
class recall	99.63%	99.16%	99.20%	97 50%	96.07%	

FIGURE 14.8: Estimated performance of the n-grams method with 10-fold cross-validation (upper image) and performance of the method on the test set (lower image).

Read more in the complete book, which can be purchased from here:

http://www.crcpress.com/product/isbn/9781482205497

...or most other book stores

Subject Index

1-grams, 211 10-Fold Cross-Validation, 220 2-grams, 213 2D Structure Descriptors, 315 3-Grams, 213 3D Scatter Plot, 162, 166, 173 3D Structure Descriptors, 315 A-B-C Segments, 8 Accuracy, 94, 275 Adjusted Rand Index, 161 Advanced Analytics, 3 Affinity, 78 Affinity-Based Marketing, 77 Agglomerative Clustering, 159 Aggregation, 84 AML, 262 AML Data Import, 262 Analogy Reasoning, 6 Analysis of Variances (ANOVA), 294 Analytics, 3 Anomaly Detection, 395 ANOVA, 294 API, 207 Area Detection, 340 Area under the Curve, 92, 94, 96 Artificial Neural Network, 290, 291 Artificial Neural Network Learner, 316 ASCII, 211 Association Rule Mining, 97, 113, 114, 234, 235, 239 Association Rule Visualization, 116 Association Rules, 22, 100, 249, 284 Astronomy, 257 Astroparticle Physics, 257 Attribute Role, 13, 14, 288 Attribute Roles, 41, 86, 321 Attribute Selection, 12, 67, 109, 150, 242, 261, 264, 322 Attribute Value Type, 40, 82, 92, 111 Attribute Value Type Transformation, 82 Attribute Value Types, 14, 321 Attribute Weighting, 22, 264, 322 Attributes, 11, 12, 14 AUC, 92, 94, 96 Audio Recommender System, 121 Automated Text Classification, 194

Backward Elimination, 323 Bag-of-Words Model, 215 **Bag-of-Words Representation**, 213 Balanced Training Set, 89 Banking Industry, 77 Bayesian Personalized Ranking Matrix Factorization, 122 Beam Search, 325 Behavior, 84 Big Data, 4 Bigrams, 213 Binary Classification, 78, 91, 95, 275 **Binomial Attribute**, 82 **Binomial Classification**, 91 **Binominal Classification**, 275 Biological Activities, 313, 314 **Biological Property**, 311 Block Plot, 173 Boolean Attribute, 82 Bootstrap Validation, 326 Bootstrapping, 22, 271, 326 Business Understanding, 78 Carcinogenicity Prediction, 314 Carpal Tunnel Syndrome (CTS), 281, 285 CART, 150 Causal Relations, 6 Centroid, 159, 160 Changing Attribute Roles, 41, 288, 321 Changing the Attribute Value Type, 40 Channel Selection, 8 Character N-Grams, 210, 223 Characteristics, 11 Chemical Properties, 314 Chemical Structures, 313 Chemistry, 319 Chemoinformatic Model, 314 Chemoinformatic Models, 311 Chemoinformatic Prediction Model, 311, 315Chemoinformatics, 311, 319 Churn Prediction, 7, 9, 94 Churn Prevention, 7, 9, 94 Class Imbalance, 88 Classification, 8, 22, 25, 33, 45, 53, 78, 145, 149, 208, 211, 214, 229, 272, 283, 288, 314, 350

456 RapidMiner: Data Mining Use Cases and Business Analytics Applications

Classification and Regression Tree (CART), 150Classification of Images, 340, 344 Classification of Text, 207 Cluster Centroids, 250 Cluster Density, 160 Cluster Internal Validation, 182 Cluster Model, 250 Cluster Validity Measures, 159 Cluster Visualization, 166, 250 Clustering, 8, 22, 157, 158, 181, 234, 240, 242, 250, 284 Clustering Validity Measures, 157 Coincidence, 4 Collaborative Filtering, 141 Collaborative Filtering Recommender System. 120 Collaborative Recommender System, 127, 130Comma Separated Values File, 101 Concept, 13 Confidence, 90, 275 Confidence Threshold, 90 Confusion Matrix, 95 Construction, 10 Content Filtering, 208 Content-Based Recommendation, 141 Content-Based Recommender System, 120. 132Contigency Table, 95 Contingency Table, 89 Conversion of Images, 338 Convert Images, 335 Correlation, 6, 22, 127, 322 Cosine Correlation, 127 Cosine Similarity, 135, 215 Cost-Based Performance Evaluation, 293 Covering Algorithm, 364 Credit Default Prediction, 8, 53 Credit Risk Scoring, 53 Credit Scoring, 8 CRISP-DM, 78 Cross-Distance, 135 Cross-Industry Standard for Data Mining, 78 Cross-Marketing, 284 Cross-Selling, 9, 284 Cross-Validation, 22, 87, 95, 151, 202, 220, 263, 275, 290, 326 CSV File, 101

CSV File Import, 46, 58, 67, 101, 166, 217, 263, 315 CTS, 281, 285 Customer Behavior, 84, 97 Customer Churn Prevention, 94 Customer Insight, 8 Customer Lifetime Value, 9 Customer Loyalty, 9, 98 Customer Profile, 12 Customer Relationship, 86 Customer Retention. 9 Customer Segmentation, 8 Customer Service Process Automation, 9 Data, 15, 16 Data Cleaning, 86, 95 Data Cleansing, 111 Data Export, 22 Data Import, 21, 39, 46, 48, 58, 67, 82, $101,\,166,\,195,\,217,\,320$ Data Import Wizard, 195 Data Loading Wizzard, 263 Data Mining, 4 Data Preparation, 81, 95, 286 Data Preprocessing, 286, 321 Data Transformation, 22, 111 Data Type, 92 Data Types, 14, 82, 92 Data Understanding, 79 Data Warehouse, 79 Database, 21–23 Database Import, 105 Dataset, 16 Davies Bouldien, 160 DBSCAN, 182 Dcoument Frequency, 208 Decay Parameter, 292 Decision Support, 284 Decision Tree, 25, 27, 150, 272, 289, 316, 323, 341, 345 Decision Tree Induction, 25, 27 Decision Trees, 92 Demand Forecasting, 10 Deployment, 93, 95, 203 Detecting Text Message Spam, 193 Diabetes, 281, 282 **Dimensionality Reduction**, 339 Direct Mail, 77 Direct Mailing, 284 Direct Marketing, 8, 284

Direct Marketing Campaign Optimization, 8, 77 Discretization, 61, 87, 145 Distance Measure, 127 Distance-Based Decision Tree, 364 Document Frequency, 236 Document Representation, 211 Document Vector, 222, 236 Document Vector Model, 213 Download, 19, 35, 45, 48, 58, 65, 126, 136-138, 140, 141, 164, 167, 195, 216, 217, 234, 261, 262, 312, 334 Drug Design, 319 Drug Effect Prediction, 320 Dummy Coding, 93 E-Coli Data, 159, 161-163, 167, 176 E-Commerce, 7 Edge Detection, 340, 350 Edge Enhancement, 340 Educational Data Mining, 143, 145, 181 Effect Coding, 93 Electronics, 10 Encoding, 211 Ensemble Classifier, 150 Ensemble of Classifiers, 272 Entropy, 28 Error Prediction, 8 ETL, 207, 228 Euclidean Distance, 215 Evaluating Feature Selection Algorithms, 264Evaluating Feature Selection Stability, 267Evaluating Feature Weighting Algorithms, 264Evaluation, 22, 37, 48, 50, 61, 69, 87, 146, 148, 150, 202, 292, 325 Example, 14, 95 Example Selection, 363 Example Set, 15, 249, 347 Example Weights, 88 Examples, 13, 15 Excel File Import, 23, 67, 315 Export, 22 Export Images, 335 Extensions, 235, 334 Fact Table, 80 Factorization-Based Recommender System, 128

Failure Prediction, 8, 10 False Negatives, 95 False Positive Rate, 95 False Positives, 95, 275 Feature Extraction, 234, 334, 338, 339, 349Feature Selection, 150, 261, 264, 322 Features, 11 Feautre Selection Stability Validation, 267Feed-Forward Backpropagation Neural Network, 291 Filter Examples, 271 Filtering Examples, 287 Finance Sector, 194 Financial Services, 7 Forward Selection, 264, 323 Fowlkes-Mallow Index, 161 FP-Growth, 113, 114, 239 Fraud Detection, 7 Frequency Distribution of Words, 250 Frequent Item Set, 113 Frequent Item Set Mining, 9, 239 Gaussian Blur, 335, 350 Gaussian Mixture Clusters, 162, 166 Generate Attributes, 269, 275 Generating Attributes, 82, 321 Glas Identification, 45 Global-Level Feature Extraction, 334, 339, 340, 344 Global-Level Features, 347 Graphical User Interface, 19, 235

Handling Missing Values, 322 Health Care Sector, 280 Hierarchical Clustering, 158 Hotel Review Analysis, 234 HSV, 333 HTML, 229 HTTP, 208, 228 Human Resources, 194 Hybrid Recommender System, 120, 135, 141 Hypothesis Test, 284, 294

GUI, 19, 235

ID Attribute, 14 Image Classification, 340, 344, 350 Image Combinations, 339 Image Conversion, 335, 338, 339 458 RapidMiner: Data Mining Use Cases and Business Analytics Applications

Image Data, 334, 339 Image Database, 335 Image Export, 335 Image Feature Extraction, 334, 338-340, 349Image Import, 335 Image Mining, 281, 333, 347, 349 image mining, 334 Image Mining Extension for RapidMiner, 333, 334 Image Segmentation, 340, 341 Image Transformation, 339 **Image Transformations**, 339 IMMI Extension, 333, 334 Import, 21, 23, 217 Import CSV Files, 195, 197 Import Data, 82, 166 Import Data from Database, 105 Import Images, 335 Indicator Attributes, 85 Indicators, 11 Influence Factors, 6, 11 Information Gain, 28, 322 Installation, 19, 139, 164, 194, 234, 235, 261, 312, 334 Instance Seletion, 363 Integration, 208 Intrusion Detection. 395 Item Recommendation, 122 Item Sets, 239 Iterating over a Set of Files, 337 Iterating over a Set of Images, 335, 336 Iteration. 148 Jaccard Index, 161, 268 Java Database Connectivity, 101 JDBC. 101 Join, 82, 135 k-Means, 158, 159, 167, 182 k-means Clustering, 242 k-Medoids, 158, 182 k-Nearest Neighbor, 33, 45, 131 k-Nearest Neighbors, 122, 137, 289, 364 k-Nearest Neighbours, 208, 214, 226 k-NN, 33, 45, 122, 131, 137, 208, 214, 226, 289, 364 Kennard-Stone Sampling, 271 Knowdledge Discovery from Textual Databases, 234

Kuncheva Index, 268

Label, 11, 13, 14, 86 Label Type Conversion, 321 Labeling, 337 Language Identification, 207, 209 Latent Features, 129 Learning Algorithm, 334 Learning Rate, 292 Leave-One-Out Validation, 326 Lift Chart, 90 Linear Regression, 93, 283, 289, 315, 317 Local Level Feature Extraction, 340 Local Maxima, 292 Local Minima, 292 Local Outlier Factor, 395 Local-Level Feature Extraction, 334, 339 Local-Level Features, 347 LOF, 395 Logging, 432 Logistic Model Tree, 150 Logistic Regression, 93, 323 Logistics, 8 Loop, 148, 263, 266 Loop Files, 217 Loop over Attributes, 322 Loop Parameters, 263, 266 Loyalty Cards, 99 M5 Prime, 317 Machine Failure Prediction, 10 Machine Failure Prevention, 10 Machine Learning Algorithm, 334 Machine Learning Research, 425 Machine Translation, 208 Macro Variables, 148, 336 Manudacturing Process Optimization, 10 Manufacturing, 10 Market Basket Analysis, 9, 97, 284 Marketing, 12 Marketing Campaign Optimization, 77 Markov Models, 214 Martketing, 8 Matrix Factorization, 122, 127, 128, 141 Maximum Relevance Minimum Redundancy Feature Selection, 261, 264

Media, 9 Medical Data Mining, 280

Meta Data, 15, 16

Meta-Learning, 425

Meta-learning, 436

Missing Value Handling, 86, 111

SUBJECT INDEX

Model, 16 Model Application, 203, 291 Model Updates, 131 Modeling, 16, 22, 25, 87, 288, 323 Molecular Descriptors, 311 Molecular Properties, 311, 320 Molecular Structure Formats, 311 Molecular Structures, 313 Momentum, 292 Monte Carlo Simulation, 269 Movie Recommender System, 121 MRMR Feature Selection, 261, 264, 265 Multi-Layer Neural Network, 291 Multi-Layer Perceptron, 291 Multiple Linear Regression Model, 315, 317Music Recommender System, 121 MySQL, 106 N-Grams, 210, 213, 223, 240, 250 Naïve Bayes, 149, 201, 214, 222, 289 Naive Bayes, 53, 65 Natural Language Processing, 208, 228 Nearest Neighbor, 33, 45 Nearest Neighbours, 208, 214, 226 Negative Example, 14 Neighborhood-Based Recommender System, 127 Network Analysis, 9 Neural Network, 92, 289-291 Neural Network Learner, 316 Neural Networks, 334 Neutrino Astronomy, 257 News Categorization, 194 News Filtering, 194 Next Best Action, 8 NLP, 208, 228 Normalization, 292 Nursery Data, 65 Object Detection, 335, 340 OLAP, 3 Online Analytical Processing, 3 Open Color Image, 336 Open Grayscale Image, 335 Operational Model, 16 Operator, 21 Opinion Mining, 8, 228 Optimization, 150

Optimize Parameters, 266

Optimizing Feature Selection and Machine Learning, 264 **Optimizing Throughput Rates**, 10 Outlier Detection, 7 Outlier Factor, 395 Over-Fitting, 150, 317, 322 PaDEL, 311, 320 PaDEL Extension for RapidMiner, 312 Parallelization, 436 Parameter Loop, 263, 266 Parameter Optimization, 93, 266 Partitional Clustering, 158 Patent Text Analysis, 10 PCA, 47, 50 Pearson Correlation, 127 Performance Evaluation, 38, 48, 50, 61, 69, 89, 265, 275, 292 Performance Measures, 124 Performance Metrics, 89, 95, 265 Permutation, 326 Personalized Recommender System, 120 Perspective, 19 Pharmaceutical Data Exploration Laboratory, 311 Pharmaceutical Industry, 280, 319 Physics, 257 Plotters, 173 PMML Extension for RapidMiner, 235 Point of Interest Detection, 340 Porter Stemmer, 135, 237 Ports, 22 Positive Example, 14 Precision, 89, 95, 125 Prediction, 275 Prediction of Carcinogenicity, 314 Predictive Accuracy, 275 Predictive Analytics, 8, 9 Predictive Maintenance, 10 Predictive Model, 25 Preventive Maintenance, 10 Price Prediction, 10 Principal Component Analysis, 47, 50 Probabilistic Classifier, 149 Process, 22, 23 Process Documents, 222, 223, 229, 235 Product Recommendation, 120 Product Recommendations, 9 Prototype Selection, 363 Prototype-Based Rules, 363 Pruning, 236

Quality Assurance, 10 Quality Optimization, 10 Quality Prediction, 10 Quantitative Structure-Activity Relationsship, 313R Console, 163 R Extension for RapidMiner, 235 R Packages, 164 R Script, 164, 169 Radial-Basis Function Kernel, 343 Rand Index, 161 Random Forest, 150, 265, 272, 323 Random Forest Learner, 316 Random Forests, 341 Ranking, 78, 121, 135, 215, 284 RapidAnalytics, 135, 138, 208, 228 RapidMiner, 19 RapidMiner Feature Selection Extension, 261RapidMiner Image Mining Extension, 349 RapidMiner IMMI Extension, 349 RapidMiner Instance Selection and Prototype-Based Rules Extension, 363 RapidMiner ISPR Extension, 363 RapidMiner PaDEL Extension, 320 RapidMiner R Extension, 163 RapidMiner Recommender Extension, 121 RapidMiner Text Processing Extension, 194RapidMiner Weka Extension, 150, 273 RapidMiner WhiBo Extension, 182 Rating, 121, 122 RBF, 343 Re-Balancing, 88 Reasoning by Analogy, 6 Recall, 89, 95 **Receiver Operating Characteristics**, 95 Recommender Performance Evaluation, 121Recommender Performance Measures, 124 Recommender System, 119, 121, 141, 143 Recommender System Web Service, 138 Recommender Systems, 9 Redundancy, 322 Redundant Attributes, 87 Regression, 22, 283, 315, 317 Regular Attribute, 14 Regular Expression, 212 Regular expressions, 430

QSAR, 313

Relative Validity Measures, 161 Relief, 322 Removing Useless Attributes, 323 Renaming Attributes, 321 Reporting Extension for RapidMiner, 235 Repository, 21, 239, 242, 249 Reputation Monitoring, 194 Retail, 7, 8, 97 RGB, 333, 347 Risk Analysis, 8 Risk Management, 8 ROC, 95 ROC Chart, 90, 91, 94 ROI Statistics, 343 Roles, 86 Rule-Based Model, 214 Running a Process, 242 Sales, 8, 12 Sampling, 88, 109, 271 SAR, 313 Saving a Process, 242 Saving Process Results, 249 Script, 438 SDF, 315 Segment-Level Feature Extraction, 334, 339, 340 Segment-Level Features, 347 Segmentation, 340 Select Attributes, 109, 145, 287 Select Examples, 271 Selecting Attributes, 82 Selecting Columns, 82 Selecting Examples, 287 Selecting Machine Learning Algorithms, 289.425 Sensor Data, 10 Sentence Tokenization, 208, 212 Sentiment Analysis, 8, 208 Series Plot, 173, 176 Similaritry, 22 Similarity Measure, 127, 335 Similarity Score, 135 Similarity-Based Content Recommendation, 134 Similarity-Based Model, 214, 226 Singular Value Decomposition, 137 SMILES, 315, 316 SMS, 193, 195, 197 Social Media Analysis, 208, 215 Spam Detection, 193, 195, 364

SUBJECT INDEX

Sparse Data Format, 122, 124 SPR, 313 SQL, 106 SQL Database, 140 Star Schema, 80 Statistical Analysis, 294 Statistical analysis, 433 Stemming, 135, 208, 237 Stopword Filter, 237 Stopword Removal, 135, 137 Stratification, 88 Stratified Sampling, 271 Structure-Activity Relationship, 313 Structure-Property Relationship, 313 Subprocess, 21, 145, 275 Sum of Squares Item Distribution, 160 Supply Chain Management, 8, 10 Support Vector Clustering, 182 Support Vector Machine, 208, 214, 289, 316, 334, 341, 343 SVM, 208, 214, 289, 316, 334, 341, 343 t-test, 434 Target Attribute, 11, 13, 14 Target Property, 313 Target Variable, 11, 13 Teacher Assistant Evaluation Data, 35 Telecommunications, 7, 9 Term Frequency, 208, 236 Term N-Grams, 240, 250 Text Analysis, 10 Text Categorization, 207, 234 Text Classification, 193, 194, 199, 207, 234 Text Clustering, 234, 240, 242, 250 Text Data, 207, 233 Text Document Filtering, 194 Text Message Spam, 193 Text Mining, 10, 135, 193, 197, 207, 233, 234Text Processing, 135, 200 Text Processing Extension for Rapid-Miner, 235 Text Representation, 211, 213 TF-IDF, 208, 214, 236 **TFIDF** Word Vector Representation, 135 Thermography, 281 Threshold, 90, 94 Time Series Analysis, 8 Time Series Forecasting, 8 Token, 212 Token Filter, 237

Token Length Filter, 137 Tokenization, 208, 212, 222, 237 Tokenization of Text Documents, 137 Tokenizing Text Documents, 135, 197 Trading Analytics, 8 Training, 96 Training Cycles, 291 Transform Cases, 222, 237 Transport, 8 Trend Analysis, 8, 10, 234 True Negative, 96 True Positive, 96 True Positive Rate, 96 True Positives, 275 Type Conversion, 93, 111, 321 Unicode, 208, 211 Unigrams, 211 Unstructured Data, 207, 233, 334, 339 Unsupervised Learning, 158 Up-Selling, 9 Update, 235 Updates, 234 URL, 229 User-Item Matrix, 141 UTF-8, 208, 211 Utility Matrix, 122 Validation, 22, 37, 48, 87, 146, 148, 150, 202, 220, 292, 318, 325 Value Type, 82, 92, 111 Value Type Conversion, 321 Value Type Transformation, 82, 111 Value Types, 14 Variables, 11, 336 Video Recommender System, 121, 126, 134View, 19 Virtual Drug Screening, 319 Visualization, 116, 157, 173 Wallace Indices, 161 Web Mining, 208 Web Mining Extension for RapidMiner, 208, 229 Web Page Language Identification, 228 Web Services, 138, 208, 228 Weight Attribute, 14 Weighted Regularized Matrix Factorization, 122 Weka, 273

462 RapidMiner: Data Mining Use Cases and Business Analytics Applications

Weka Extension for RapidMiner, 235
Word Frequency, 197
Word List, 135, 197, 198, 229, 249, 250
Word Stemming, 208, 237
Word Vector, 135, 197, 198, 222
Wrapper Validation, 150
Wrapper X-Validation, 263–265

X-Validation, 22, 87, 95, 151, 202, 220, 263, 275, 290, 326 XML, 208, 235

Y-Randomization, 326

Operator Index

Add Hydrogen - PaDEL, 313 Add Noise, 375 Aggregate, 295, 371 Aggregation, 84 All-knn-selection, 368 ANOVA, 294, 435 Append, 219, 263, 295, 375, 404 Apply Model, 38, 89, 131, 202, 222, 265, 274, 291, 316, 415 Apply Threshold, 90 approximate Local Correlation Integral, 408Attribute Based Item Recommendation, 136Attributes by Weight, 323 Backward Elimination, 323 Bootstrap Validation, 326 Border / Interior Classification, 353 Calculate Descriptors - PaDEL, 313 Class assigner, 381 Class Iterator, 386 Clear Log, 375 Cluster Distance Performance, 416 Cluster Internal Evaluation, 182 Cluster-Based Local Outlier Factor (CBLOF), 409Clustering (k-means(fast)), 242 Color Component Extractor, 337 Color to gray-scale, 337, 352 Combine Documents, 219 Compare ROCs, 93 Connectivity-Based Outlier Factor (COF), 407Convert to 3D - PaDEL, 313 Count Substructure - PaDEL, 313 Create Association Rules, 114, 239 Create Lift Chart, 90 Create Threshold, 90 Cross Distance, 135 DBScan, 185 Decision Tree, 25, 92, 150, 301, 429 Declare Missing Values, 111 Detect Aromaticity - PaDEL, 313

Detect Dimensionality - PaDEL, 313

Detect Outlier (Densities), 397

Detect Outlier (Distance), 397 Detector, 357 Dichomatize, 321 Discretize, 145 Discretize by Binning, 53, 56, 87 Discretize by User Specification, 87 dLog Distance, 351 Documents to Data, 219 ENN selection, 368 ENN Weighting, 384 Execute Process, 164 Execute Script, 416, 438 Extract Cluster Prototypes, 390 Extract Content, 229 Extract Macro, 275, 325, 381, 419, 438 FCM, 381 Feature Selection Stability Validation, 268Filter Examples, 53, 56, 83, 111, 147, 276, 286, 323, 343, 412 Filter Stopwords, 201 Filter Stopwords (English), 237 Filter Tokens (by Length), 237 Forward Selection, 323 FP-Growth, 113, 239 Free Memory, 265, 325 Fuzzy C-Means, 381 Gaussian Blur, 335, 358 Generate Attributes, 82, 217, 269, 275, 295, 322, 412 Generate Data, 162, 403 Generate ID, 326 Generate Macro, 381 Generate n-Grams (Characters), 224 Generate N-Grams (Terms), 200, 240 Generate Weight (Stratification), 88 Get Page, 229 Global Feature Extractor from a Single Image, 345, 351 Global statistics, 345, 352 Grayscale to color, 337 Grouped ANOVA, 303 Guess Type, 85

Handle Exception, 375, 435

Histogram, 352 Histogram Equalization, 336

IB3 Selection, 370 Image Combinator, 339 Impute Missing Values, 87 Influenced Outlierness (INFLO), 408 Interest Points Visualizer, 357 Inverse, 343 Item Attribute k-NN, 137 Item k-NN, 131, 137 Iterative Component Analysis, 47

Join, 82, 135, 323, 413

K-Means, 167, 185, 300, 417
K-medoids, 185
k-NN, 33, 42, 45, 371, 381, 431
k-NN Global Anomaly Score, 404
Keep Document Parts, 212

LDA, 375 Learning Vector Quantization, 381 LibSVM, 414 Linear Discriminant Analysis, 375 Linear Regression, 93, 317 Local Correlation Integral (LOCI), 408 Local Density Cluster-Based Outlier Factor (LDCOF), 410Local Outlier Factor (LOF), 406 Local Outlier Probability (LoOP), 408 Log, 89, 148, 265, 325, 371, 432 Log to Data, 138, 325, 371, 417, 433 Logistic Regression, 93, 323 Look Up table applier, 337 Loop, 148, 267, 336, 381, 431 Loop and Deliver Best, 323 Loop Attributes, 322 Loop Datasets, 375 Loop Examples, 417 Loop Files, 217, 336, 429 Loop Labels, 384 Loop Models, 387 Loop Parameters, 185, 263, 266, 375, 387 Loop Subset, 325 LVQ, 381 Macro, 336

Map, 381, 403, 415 Materialize Data, 167, 387 MC Selection, 367 Merge Segments to Image, 343 Model Combiner, 136 Modelling, 315 Multiple Color Image Opener, 345, 350 Multiple Grayscale Image Opener, 336 Multiply, 83, 135, 147, 271, 295, 323, 338, 413Naive Bayes, 53, 65, 150, 202, 222, 431 Neural Net, 290, 291 Nominal to Binominal, 321 Nominal to Binominal, 82, 111 Nominal to Numerical, 93 Nominal to Text, 219 Normalize, 185, 371, 399, 413 Numerical to Binominal, 37, 40, 237, 295 Numerical to Nominal, 381 Numerical to Polynomial, 37, 40 Open Color Image, 336 Open Grayscale Image, 335, 357 Optimize Parameters, 93, 187 Optimize Parameters (Grid), 415 Optimize Search, 325 Optimize Selection, 150, 325 Order-based Block Color, 352 PaDEL, 312 Parse Numbers, 138, 378 Performance, 38, 51, 89, 131, 137, 202, 265, 292, 429 Performance (Binominal Classification), 53, 89 Performance (Classification), 38, 48, 148, 415Performance - Regression, 317 Performance Binominal Clasification, 38 Permutation, 326 Pivot, 433 poi_generator, 346 Principal Component Analysis, 47 Process Documents, 135, 197 Process Documents from Data, 135, 197, 221Process Documents from Files, 235 Provide Macro as Log Value, 138, 381, 419, 432 Random Forest, 150, 273, 323, 431 Random Selection, 367, 384 Read AML, 124, 145, 183, 262, 428 Read Compounds - PaDEL, 313 Read CSV, 36, 46, 105, 166, 197, 217, 411

OPERATOR INDEX

Read Database, 82, 105 Read Excel, 23, 36, 67 Read Image Set, 356 Read URL, 36 Recall, 375, 416 Remap Binominals, 270, 321 Remember, 167, 375, 416 Remove Correlated Attributes, 47 Remove Salt - PaDEL, 313 Remove Unused Values, 286 Remove Useless Attributes, 47, 87, 323 Rename, 36, 40, 166, 217, 262, 295, 321 Rename by Replacing, 53, 56 Rename Compound - PaDEL, 313 **RENN** selection, 375 Replace, 112 Replace Missing Values, 86, 111, 185, 299 Replace Tokens, 224 Replenish Values, 322 Retrieve, 82, 197, 220, 321 RMHC selection, 367 RNG selection, 370 **ROI** Statistics, 358 ROI statistics, 343 Sample, 88, 109, 200, 271 Sample (Bootstrapping), 271 Sample (Stratified), 387 Segment Feature Extraction, 343 Segment Feature Extractor, 358 Segment Filter by Example Set, 343 Segment Mask Processor, 358 Select, 434 Select Attributes, 23, 47, 67, 82, 109, 145, 167, 183, 217, 242, 286, 411 Select by MRMR/CFS, 265 Select by Weights, 323 Select Subprocess, 148, 185, 378, 431 Set Minus, 135 Set Role, 36, 37, 86, 130, 145, 217, 219, 326, 412 Simple Validation, 148 Singular Value Decomposition, 47 Sort, 322 Split Data, 48, 219 Split Validation, 37 Statistical region merging, 343, 358 Stem (Porter), 200, 237 Store, 197, 219, 322, 343, 404 Subprocess, 21, 145, 275, 343, 432 Support Vector Clustering, 185

Support Vector Machine, 431 Support Vector Machine (LibSVM), 403, 415SVD, 137 SVM, 342, 387 T-Test, 434 Text to Nominal, 237 Thresholding, 358 Tokenize, 197, 220, 237 Trainable segmentation, 342 Transform Cases, 200, 226, 237 Transpose, 323 Unescape HTML, 229 Update Model, 131 Validation, 221 Vector Quantization, 381 Viola-Jones, 357 VQ, 381 W-LMT, 150 W-Random Forest, 265 W-SimpleCart, 150 Weight by Correlation, 323 Weight by Information Gain, 323 Weight by Tree Importance, 323 Weight Transformation, 386 Wrapper Split Validation, 150 Wrapper X-Validation, 265 Write AML, 276 Write Compounds - PaDEL, 313 Write Image to File, 335, 336 Write Model, 94 Write Weights, 151, 263 X-Means, 409 X-Validation, 38, 53, 57, 87, 203, 222, 263, 290, 429