

# Table of Contents

<b>FOREWORD by William Hersh</b> . . . . .	<b>.xi</b>	<b>CHAPTER 2: DATA PREPARATION AND WRANGLING</b> . . . . .	<b>23</b>
<b>PREFACE</b> . . . . .	<b>.xiii</b>	Understanding Data And Data Types. . . . .	24
<b>CHAPTER 1: INTRODUCTION TO DATA SCIENCE</b> . . . . .	<b>1</b>	Introduction . . . . .	24
What Is Data Science? . . . . .	1	Data Science DISCOVER Process Step 2 . . . . .	24
Background . . . . .	2	Data Types . . . . .	27
Artificial Intelligence and Machine Learning. . . . .	2	Step 3 In The Data Science Process . . . . .	29
Artificial Intelligence and Data Science . . . . .	5	Data Cleaning . . . . .	29
Foundations and Frontiers: Integrating DMBok Practices with Data Science. . . . .	5	What is Feature Engineering? . . . . .	29
Statistical Foundations and Predictive Frontiers: A Data Voyage . . . . .	6	Handling Missing Data with Imputation or Deletion . . . . .	30
The 8-Step Data Science Discover Process . . . . .	6	Handling Outliers . . . . .	33
The 3 Levels of AI for Healthcare . . . . .	7	Binning / Discretize . . . . .	37
What Do Data Scientists Do? . . . . .	9	Log Transform . . . . .	41
Why Learn Data Science? . . . . .	10	Feature Creation . . . . .	42
Data Science Educational Challenges . . . . .	10	One Hot Encoding Versus N-1 Dummy Encoding . . . . .	42
Why Did We Create This Textbook? . . . . .	11	Group Calculation Operations . . . . .	45
Who Is The Target Audience? . . . . .	12	Feature Splitting . . . . .	45
Why Use Open-Source Software? . . . . .	12	Feature Scaling. . . . .	45
Orange Data Mining . . . . .	13	Univariate, Bivariate, And Multivariate Analysis . . . . .	49
Orange General Information. . . . .	13	Important Data Challenges . . . . .	50
JASP Stats Package . . . . .	16	Curse of Dimensionality . . . . .	50
Introduction . . . . .	16	Multicollinearity . . . . .	50
JASP Unique Features . . . . .	16	Data Leakage . . . . .	62
JASP Organization. . . . .	17	Imbalanced Data . . . . .	63
JASP limitations . . . . .	18	Synthetic Data . . . . .	63
Orange vs JASP . . . . .	18	Introduction . . . . .	63
Putting Knowledge Into Practice . . . . .	18	Synthea. . . . .	63
Exercise . . . . .	18	AI-Generated Synthetic Tabular Data . . . . .	64
Maturity Levels . . . . .	19	Data Science Experiments (DSEs) . . . . .	66
Resources . . . . .	19	Experimental Thinking in Data Science. . . . .	66
References . . . . .	20	Nine Steps for Creating Effective Data Science Experiments . . . . .	67
		Data Science Experiment Complexity Level . . . . .	68

Algorithm Preprocessing Defaults . . . . .	67	Types of Gradient Descent . . . . .	151
Algorithm Compatibility with Other Widgets . . . . .	69	Orange Algorithms . . . . .	151
Data Science Experiment Complexity Level Examples . . . . .	70	JASP Models . . . . .	159
Orange Data And Transform Menu Widgets . . . . .	77	Linear Regression . . . . .	160
JASP Data Features . . . . .	97	Logistic Regression . . . . .	160
Putting Knowledge Into Practice . . . . .	97	JASP Machine Learning Algorithms . . . . .	160
Exercise . . . . .	97	MLOps . . . . .	161
Maturity Level . . . . .	98	Tackling Model Challenges . . . . .	162
Resources . . . . .	99	Model Bias and Fairness . . . . .	162
References . . . . .	99	Navigating Imbalanced Datasets . . . . .	162
		Hyperparameter Optimization . . . . .	173
		Avoiding Model Overfitting . . . . .	173
		Putting Knowledge Into Practice . . . . .	173
		Exercise . . . . .	173
		Maturity Levels . . . . .	175
		Resources . . . . .	175
		References . . . . .	176
<b>CHAPTER 3: DATA VISUALIZATION . . . . .</b>	<b>105</b>		
Introduction . . . . .	106	<b>CHAPTER 5: MODEL EVALUATION . . . . .</b>	<b>179</b>
Data Visualization Tables And Plots . . . . .	109	Model Validation . . . . .	180
Statistical Tables and Displays . . . . .	109	Internal Validation (IV) . . . . .	180
Basic Data Displays . . . . .	110	External Validation (EV) . . . . .	181
Distributions . . . . .	113	Optimal Validation (OV) . . . . .	182
Relationships Between Variables . . . . .	117	Further Validation Guidance . . . . .	182
Time Series . . . . .	122	Model Discrimination And Calibration . . . . .	182
Geographical Displays . . . . .	122	Probability, Odds, Odds Ratios And	
Hierarchical Displays . . . . .	124	Log Odds Ratios . . . . .	185
Multivariate Displays . . . . .	124	Probability vs Odds . . . . .	185
Statistical Models and Displays . . . . .	129	Odds Ratios and Log Odds Ratios . . . . .	185
Text Displays . . . . .	135	Diagnostic Tests and Likelihood Ratios . . . . .	187
Variable Ranks - Non-Regularized Versus Regularized . . . . .	136	Orange Model Evaluation Widgets . . . . .	188
Dashboards And Infographics . . . . .	136	Test and Score Widget . . . . .	189
Dashboards . . . . .	136	Predictions Widget . . . . .	191
Infographics . . . . .	138	Confusion Matrix . . . . .	192
Putting Knowledge Into Practice . . . . .	140	ROC Curve Widget . . . . .	193
Exercise . . . . .	140	Precision-Recall (PR) Curve . . . . .	196
Maturity Levels . . . . .	141	Other Valuable Orange Curves . . . . .	196
Resources . . . . .	141	Lift Curve . . . . .	196
References . . . . .	141	Calibration Plot . . . . .	196
		Residual Plots . . . . .	198
		Explain Module . . . . .	201
		Rank Widget . . . . .	201
		Feature Importance Widget . . . . .	203
		Explain Prediction Widget . . . . .	204
		Explain Model Widget . . . . .	206
<b>CHAPTER 4: MACHINE LEARNING MODELS . . . . .</b>	<b>143</b>		
Introduction . . . . .	143		
How Do Models Work? . . . . .	146		
More Modeling Concepts . . . . .	147		
Bias-Variance Tradeoff . . . . .	147		
Model Interpretability and Explainability . . . . .	149		
Gradient Descent . . . . .	150		
Local and Global Minimum . . . . .	150		

JASP Model Evaluation . . . . .	206	PCA Exercises . . . . .	239
Hypothesis Testing . . . . .	208	Association Rules . . . . .	241
Hypothesis Testing Example . . . . .	208	Association Rules Exercise . . . . .	242
How Hypothesis Tests Make Decisions . . . . .	208	Putting Knowledge Into Practice . . . . .	243
Type 1 and Type 2 Hypothesis Testing Errors . . . . .	208	Exercise . . . . .	243
Hypothesis Tests Available in Orange, JASP, and BlueSky Statistics . . . . .	209	Maturity Levels . . . . .	246
Identifying And Addressing Concept Drift And Model Drift . . . . .	212	Resources . . . . .	247
Concept Drift . . . . .	212	References . . . . .	247
Model Drift . . . . .	212		
When to Retrain a Model . . . . .	213	<b>CHAPTER 8: TIME SERIES FORECASTING AND SURVIVAL ANALYSIS . . . . . 249</b>	
Putting Knowledge Into Practice . . . . .	213	Time Series Forecasting With Arima . . . . .	250
Exercise . . . . .	213	Arima Time Series Forecasting in Orange . . . . .	250
Maturity Levels . . . . .	214	Pros and Cons of ARIMA Models . . . . .	252
Resources . . . . .	215	Survival Analysis . . . . .	254
References . . . . .	215	The General Steps to Perform Kaplan-Meier Survival Analysis . . . . .	254
		Orange Example of Kaplan-Meier Survival Analysis . . . . .	255
<b>CHAPTER 6: SUPERVISED LEARNING . . . . 217</b>		Putting Knowledge Into Practice . . . . .	259
Classification Modeling . . . . .	218	Exercise . . . . .	259
Introduction . . . . .	218	Maturity Levels . . . . .	260
Regression Modeling . . . . .	219	Resources . . . . .	261
Linear regression . . . . .	219	References . . . . .	261
Simple Linear Regression (SLR) . . . . .	219		
Multiple Linear Regression (MLR) . . . . .	221	<b>CHAPTER 9: GEOLOCATION . . . . . 262</b>	
The Dummy Variable Trap . . . . .	222	Introduction . . . . .	264
Logistic Regression . . . . .	223	Geolocation Exercises . . . . .	265
Polynomial Regression . . . . .	223	Geolocation with Synthetic Data . . . . .	275
Linear Regression Model Regularization . . . . .	224	Multi-Industry Application Examples Of Geo-Mapping And Geo-Analysis . . . . .	267
Regression Performance Metrics . . . . .	224	Putting Knowledge Into Practice . . . . .	269
Putting Knowledge Into Practice . . . . .	225	Exercise . . . . .	269
Exercises . . . . .	225	Maturity Levels . . . . .	270
Maturity Levels . . . . .	228	Resources . . . . .	271
Resources . . . . .	229	References . . . . .	271
References . . . . .	229		
		<b>CHAPTER 10: IMAGE ANALYTICS . . . . . 273</b>	
<b>CHAPTER 7: UNSUPERVISED LEARNING . . 231</b>		Introduction . . . . .	273
Unsupervised Learning . . . . .	232	Medical Field Applications . . . . .	273
K-Means Clustering . . . . .	232	Non- Medical Field Applications . . . . .	274
K-Means Clustering Exercises . . . . .	233	AI Imaging: Industry Application Examples . . . . .	274
Hierarchical Clustering . . . . .	237	How Do CNNs Work? . . . . .	275
Principle Component Analysis (PCA) . . . . .	237	The Normal Steps Of Image Analytics With a CNN . . . . .	275
Introduction . . . . .	237		
Routine Steps to Achieve PCA . . . . .	238		

CNN Layers . . . . .	276
Putting Knowledge Into Practice . . . . .	277
Exercises . . . . .	277
Maturity Levels . . . . .	280
Resources . . . . .	281
References . . . . .	281
<b>CHAPTER 11: TEXT MINING . . . . .</b>	<b>283</b>
Introduction . . . . .	284
Neural Networks And NLP . . . . .	285
Recurrent Neural Networks . . . . .	285
Transformer and Attention Models . . . . .	286
NLP Concepts . . . . .	286
NLP Preprocessing Steps . . . . .	287
Topic Analysis . . . . .	288
Large Language Models (LLMs) . . . . .	288
Graphs From Text . . . . .	289
Code Assistants . . . . .	291
Small Language Models (SLM) . . . . .	291
Orange NLP Widgets . . . . .	291
Text Document Widgets . . . . .	292
Preprocess Text Widget . . . . .	292
Document Embedding Widget . . . . .	292
Bag of Words Widget . . . . .	292
Topic Modeling Widget . . . . .	292
Putting Knowledge Into Practice . . . . .	292
Exercises . . . . .	292
Maturity Levels . . . . .	298
Resources . . . . .	299
References . . . . .	299
<b>CHAPTER 12: INTEGRATING CONTINUOUS IMPROVEMENT AND DATA SCIENCE INTO INDUSTRY 4.0 . . . . .</b>	<b>301</b>
Introduction . . . . .	302

Industry And Sector 4.0 Initiatives . . . . .	303
Accelerating Industry 4.0: The Pivotal Role Of Continuous Improvement (CI) Programs . . . . .	305
Eleven Powerful Continuous Improvement Techniques . . . . .	307
Hoshin Kanri with Catchball Strategic Planning . . . . .	307
Six Sigma . . . . .	308
Lean / Kaizen . . . . .	311
Lean Six Sigma . . . . .	314
Predictive Modeling . . . . .	315
High Reliability Organizations (HROs) . . . . .	316
Agile Versus Waterfall Thinking . . . . .	318
Cascading Risk Management (CRM) . . . . .	319
AI Systems . . . . .	321
AI Chatbots . . . . .	322
Innovation-On-Demand Techniques . . . . .	324
The Redwood Approach: Building Business Resilience Through Cooperative Strength. . . . .	327
Maturity Levels . . . . .	327
Resources . . . . .	328
References . . . . .	329

**APPENDIX A: DATA VISUALIZATION . . . . . 335**

**APPENDIX B: MACHINE LEARNING  
ALGORITHMS . . . . . 337**

**APPENDIX C: GLOSSARY . . . . . 341**

**INDEX . . . . . 347**