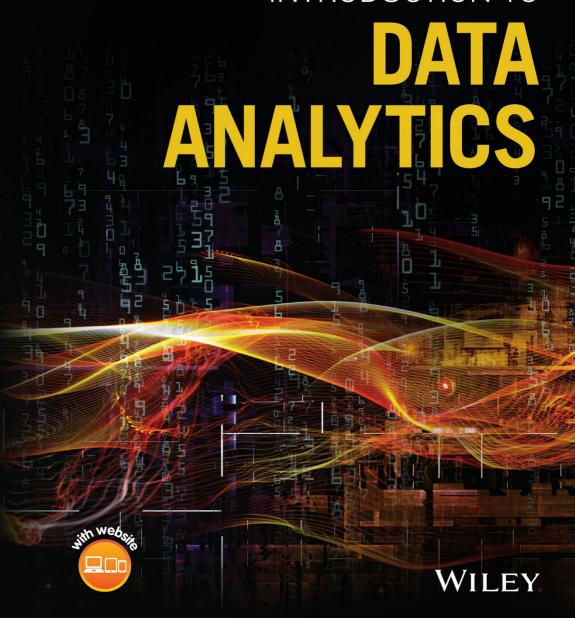
João Moreira André de Carvalho Tomáš Horváth

A GENERAL INTRODUCTION TO





A General Introduction to Data Analytics

João Mendes Moreira University of Porto

André C. P. L. F. de Carvalho University of São Paulo

Tomáš Horváth Eötvös Loránd University in Budapest Pavol Jozef Šafárik University in Košice



This edition first published 2019 © 2019 John Wiley & Sons, Inc.

All rights reserved. No part of this publication may be reproduced, stored in a retrieval system, or transmitted, in any form or by any means, electronic, mechanical, photocopying, recording or otherwise, except as permitted by law. Advice on how to obtain permission to reuse material from this title is available at http://www.wiley.com/go/permissions.

The right of João Moreira, André de Carvalho, and Tomáš Horváth to be identified as the author(s) of this work has been asserted in accordance with law.

Registered Office

John Wiley & Sons, Inc., 111 River Street, Hoboken, NJ 07030, USA

Editorial Office

111 River Street, Hoboken, NJ 07030, USA

For details of our global editorial offices, customer services, and more information about Wiley products visit us at www.wiley.com.

Wiley also publishes its books in a variety of electronic formats and by print-on-demand. Some content that appears in standard print versions of this book may not be available in other formats.

Limit of Liability/Disclaimer of Warranty

In view of ongoing research, equipment modifications, changes in governmental regulations, and the constant flow of information relating to the use of experimental reagents, equipment, and devices, the reader is urged to review and evaluate the information provided in the package insert or instructions for each chemical, piece of equipment, reagent, or device for, among other things, any changes in the instructions or indication of usage and for added warnings and precautions. While the publisher and authors have used their best efforts in preparing this work, they make no representations or warranties with respect to the accuracy or completeness of the contents of this work and specifically disclaim all warranties, including without limitation any implied warranties of merchantability or fitness for a particular purpose. No warranty may be created or extended by sales representatives, written sales materials or promotional statements for this work. The fact that an organization, website, or product is referred to in this work as a citation and/or potential source of further information does not mean that the publisher and authors endorse the information or services the organization, website, or product may provide or recommendations it may make. This work is sold with the understanding that the publisher is not engaged in rendering professional services. The advice and strategies contained herein may not be suitable for your situation. You should consult with a specialist where appropriate. Further, readers should be aware that websites listed in this work may have changed or disappeared between when this work was written and when it is read. Neither the publisher nor authors shall be liable for any loss of profit or any other commercial damages, including but not limited to special, incidental, consequential, or other damages.

Library of Congress Cataloging-in-Publication Data

Names: Moreira, João, 1969– author. | Carvalho, André Carlos Ponce de Leon Ferreira, author. | Horváth, Tomáš, 1976– author.

Title: A general introduction to data analytics / by João Mendes Moreira, André C. P. L. F. de Carvalho,

Description: Hoboken, NJ: John Wiley & Sons, 2019. | Includes bibliographical references and index. | Identifiers: LCCN 2017060728 (print) | LCCN 2018005929 (ebook) | ISBN

9781119296256 (pdf) | ISBN 9781119296263 (epub) | ISBN 9781119296249 (cloth)

Subjects: LCSH: Mathematical statistics—Methodology. | Electronic data processing. | Data mining. Classification: LCC QA276.4 (ebook) | LCC QA276.4 .M664 2018 (print) | DDC 519.50285–dc23 LC record available at https://lccn.loc.gov/2017060728

Cover image: © agsandrew/Shutterstock Cover design by Wiley

Printed in the United States of America.

Set in 10/12pt Warnock by SPi Global, Pondicherry, India

To the women at home that make my life better: Mamã, Yá and Yé – João To my family, Valeria, Beatriz, Gabriela and Mariana – André To my wife Danielle – Tomáš

Contents

Preface xiii

Acknowledgments xv

	Presentational Conventions $xvii$ About the Companion Website xix	
	Part I Introductory Background 1	
1	What Can We Do With Data? 3	
1.1	Big Data and Data Science 4	
1.2	Big Data Architectures 5	
1.3	Small Data 6	
1.4	What is Data? 7	
1.5	A Short Taxonomy of Data Analytics 9	
1.6	Examples of Data Use 10	
1.6.1	Breast Cancer in Wisconsin 11	
1.6.2	Polish Company Insolvency Data 11	
1.7	A Project on Data Analytics 12	
1.7.1	A Little History on Methodologies for Data Analytics	12
1.7.2	The KDD Process 14	
1.7.3	The CRISP-DM Methodology 15	
1.8	How this Book is Organized 16	
1.9	Who Should Read this Book 18	
	Part II Getting Insights from Data 19	
2	Descriptive Statistics 21	
2.1	Scale Types 22	
2.2	Descriptive Univariate Analysis 25	

Univariate Frequencies 25

2.2.1

viii	Contents
viii	Contents

2.2.2 2.2.3 2.2.4 2.3 2.3.1 2.3.2 2.3.3 2.4 2.5	Univariate Data Visualization 27 Univariate Statistics 32 Common Univariate Probability Distributions 38 Descriptive Bivariate Analysis 40 Two Quantitative Attributes 41 Two Qualitative Attributes, at Least one of them Nominal 45 Two Ordinal Attributes 46 Final Remarks 47 Exercises 47
3 3.1 3.2	Descriptive Multivariate Analysis 49 Multivariate Frequencies 49 Multivariate Data Visualization 50
3.3	Multivariate Statistics 59
3.3.1	Location Multivariate Statistics 59
3.3.2	Dispersion Multivariate Statistics 60
3.4	Infographics and Word Clouds 66
3.4.1	Infographics 66
3.4.2	Word Clouds 67
3.5	Final Remarks 67
3.6	Exercises 68
_	D . O . U ID
4	Data Quality and Preprocessing 71
4 4.1	• , ,
4.1	Data Quality 71
4.1 4.1.1	Data Quality 71 Missing Values 72
4.1 4.1.1 4.1.2	Data Quality 71 Missing Values 72 Redundant Data 74
4.1 4.1.1 4.1.2 4.1.3	Data Quality 71 Missing Values 72 Redundant Data 74 Inconsistent Data 75 Noisy Data 76 Outliers 77
4.1 4.1.1 4.1.2 4.1.3 4.1.4	Data Quality 71 Missing Values 72 Redundant Data 74 Inconsistent Data 75 Noisy Data 76 Outliers 77 Converting to a Different Scale Type 77
4.1 4.1.1 4.1.2 4.1.3 4.1.4 4.1.5	Data Quality 71 Missing Values 72 Redundant Data 74 Inconsistent Data 75 Noisy Data 76 Outliers 77 Converting to a Different Scale Type 77 Converting Nominal to Relative 78
4.1 4.1.1 4.1.2 4.1.3 4.1.4 4.1.5 4.2 4.2.1 4.2.2	Data Quality 71 Missing Values 72 Redundant Data 74 Inconsistent Data 75 Noisy Data 76 Outliers 77 Converting to a Different Scale Type 77 Converting Nominal to Relative 78 Converting Ordinal to Relative or Absolute 81
4.1 4.1.1 4.1.2 4.1.3 4.1.4 4.1.5 4.2 4.2.1 4.2.2 4.2.3	Data Quality 71 Missing Values 72 Redundant Data 74 Inconsistent Data 75 Noisy Data 76 Outliers 77 Converting to a Different Scale Type 77 Converting Nominal to Relative 78 Converting Ordinal to Relative or Absolute 81 Converting Relative or Absolute to Ordinal or Nominal 82
4.1 4.1.1 4.1.2 4.1.3 4.1.4 4.1.5 4.2 4.2.1 4.2.2 4.2.3 4.3	Data Quality 71 Missing Values 72 Redundant Data 74 Inconsistent Data 75 Noisy Data 76 Outliers 77 Converting to a Different Scale Type 77 Converting Nominal to Relative 78 Converting Ordinal to Relative or Absolute 81 Converting Relative or Absolute to Ordinal or Nominal 82 Converting to a Different Scale 83
4.1 4.1.1 4.1.2 4.1.3 4.1.4 4.1.5 4.2 4.2.1 4.2.2 4.2.3 4.3 4.4	Data Quality 71 Missing Values 72 Redundant Data 74 Inconsistent Data 75 Noisy Data 76 Outliers 77 Converting to a Different Scale Type 77 Converting Nominal to Relative 78 Converting Ordinal to Relative or Absolute 81 Converting Relative or Absolute to Ordinal or Nominal 82 Converting to a Different Scale 83 Data Transformation 85
4.1 4.1.1 4.1.2 4.1.3 4.1.4 4.1.5 4.2 4.2.1 4.2.2 4.2.3 4.3 4.4 4.5	Data Quality 71 Missing Values 72 Redundant Data 74 Inconsistent Data 75 Noisy Data 76 Outliers 77 Converting to a Different Scale Type 77 Converting Nominal to Relative 78 Converting Ordinal to Relative or Absolute 81 Converting Relative or Absolute to Ordinal or Nominal 82 Converting to a Different Scale 83 Data Transformation 85 Dimensionality Reduction 86
4.1 4.1.1 4.1.2 4.1.3 4.1.4 4.1.5 4.2 4.2.1 4.2.2 4.2.3 4.3 4.4 4.5 4.5.1	Data Quality 71 Missing Values 72 Redundant Data 74 Inconsistent Data 75 Noisy Data 76 Outliers 77 Converting to a Different Scale Type 77 Converting Nominal to Relative 78 Converting Ordinal to Relative or Absolute 81 Converting Relative or Absolute to Ordinal or Nominal 82 Converting to a Different Scale 83 Data Transformation 85 Dimensionality Reduction 86 Attribute Aggregation 88
4.1 4.1.1 4.1.2 4.1.3 4.1.4 4.1.5 4.2 4.2.1 4.2.2 4.2.3 4.3 4.4 4.5 4.5.1	Data Quality 71 Missing Values 72 Redundant Data 74 Inconsistent Data 75 Noisy Data 76 Outliers 77 Converting to a Different Scale Type 77 Converting Nominal to Relative 78 Converting Ordinal to Relative or Absolute 81 Converting Relative or Absolute to Ordinal or Nominal 82 Converting to a Different Scale 83 Data Transformation 85 Dimensionality Reduction 86 Attribute Aggregation 88 Principal Component Analysis 88
4.1 4.1.1 4.1.2 4.1.3 4.1.4 4.1.5 4.2 4.2.1 4.2.2 4.2.3 4.3 4.4 4.5 4.5.1 4.5.1.1	Data Quality 71 Missing Values 72 Redundant Data 74 Inconsistent Data 75 Noisy Data 76 Outliers 77 Converting to a Different Scale Type 77 Converting Nominal to Relative 78 Converting Ordinal to Relative or Absolute 81 Converting Relative or Absolute to Ordinal or Nominal 82 Converting to a Different Scale 83 Data Transformation 85 Dimensionality Reduction 86 Attribute Aggregation 88 Principal Component Analysis 88 Independent Component Analysis 91
4.1 4.1.1 4.1.2 4.1.3 4.1.4 4.1.5 4.2 4.2.1 4.2.2 4.2.3 4.3 4.4 4.5 4.5.1 4.5.1.1 4.5.1.2 4.5.1.3	Data Quality 71 Missing Values 72 Redundant Data 74 Inconsistent Data 75 Noisy Data 76 Outliers 77 Converting to a Different Scale Type 77 Converting Nominal to Relative 78 Converting Ordinal to Relative or Absolute 81 Converting Relative or Absolute to Ordinal or Nominal 82 Converting to a Different Scale 83 Data Transformation 85 Dimensionality Reduction 86 Attribute Aggregation 88 Principal Component Analysis 88 Independent Component Analysis 91 Multidimensional Scaling 91
4.1 4.1.1 4.1.2 4.1.3 4.1.4 4.1.5 4.2 4.2.1 4.2.2 4.2.3 4.3 4.4 4.5 4.5.1 4.5.1.1 4.5.1.2 4.5.1.3 4.5.2	Data Quality 71 Missing Values 72 Redundant Data 74 Inconsistent Data 75 Noisy Data 76 Outliers 77 Converting to a Different Scale Type 77 Converting Nominal to Relative 78 Converting Ordinal to Relative or Absolute 81 Converting Relative or Absolute to Ordinal or Nominal 82 Converting to a Different Scale 83 Data Transformation 85 Dimensionality Reduction 86 Attribute Aggregation 88 Principal Component Analysis 88 Independent Component Analysis 91 Multidimensional Scaling 91 Attribute Selection 92
4.1 4.1.1 4.1.2 4.1.3 4.1.4 4.1.5 4.2 4.2.1 4.2.2 4.2.3 4.3 4.4 4.5 4.5.1 4.5.1.1 4.5.1.2 4.5.1.3	Data Quality 71 Missing Values 72 Redundant Data 74 Inconsistent Data 75 Noisy Data 76 Outliers 77 Converting to a Different Scale Type 77 Converting Nominal to Relative 78 Converting Ordinal to Relative or Absolute 81 Converting Relative or Absolute to Ordinal or Nominal 82 Converting to a Different Scale 83 Data Transformation 85 Dimensionality Reduction 86 Attribute Aggregation 88 Principal Component Analysis 88 Independent Component Analysis 91 Multidimensional Scaling 91

4.5.2.4	Search Strategies 95
4.6	Final Remarks 96
4.7	Exercises 96
5	Clustering 99
5.1	Distance Measures 100
5.1.1	Differences between Values of Common Attribute Types 101
5.1.2	Distance Measures for Objects with Quantitative Attributes 103
5.1.3	Distance Measures for Non-conventional Attributes 104
5.2	Clustering Validation 107
5.3	Clustering Techniques 108
5.3.1	K-means 110
5.3.1.1	Centroids and Distance Measures 110
5.3.1.2	How K-means Works 111
5.3.2	DBSCAN 115
5.3.3	Agglomerative Hierarchical Clustering Technique 117
5.3.3.1	Linkage Criterion 119
5.3.3.2	Dendrograms 120
5.4	Final Remarks 122
5.5	Exercises 123
6	Frequent Pattern Mining 125
6.1	Frequent Itemsets 127
6.1.1	Setting the <i>min_sup</i> Threshold 128
6.1.2	Apriori – a Join-based Method 131
6.1.3	Eclat 133
6.1.4	FP-Growth 134
6.1.5	Maximal and Closed Frequent Itemsets 138
6.2	Association Rules 139
6.3	Behind Support and Confidence 142
6.3.1	Cross-support Patterns 143
6.3.2	Lift 144
6.3.3	Simpson's Paradox 145
6.4	Other Types of Pattern 147
6.4.1	Sequential patterns 147
6.4.2	Frequent Sequence Mining 148
6.4.3	Closed and Maximal Sequences 148
6.5	Final Remarks 149
6.6	Exercises 149
7	Cheat Sheet and Project on Descriptive Analytics 151
7.1	Cheat Sheet of Descriptive Analytics 151
711	On Data Summarization 151

7.1.2	On Clustering 151
7.1.3	On Frequent Pattern Mining 153
7.2	Project on Descriptive Analytics 154
7.2.1	Business Understanding 154
7.2.2	Data Understanding 155
7.2.3	Data Preparation 155
7.2.4	Modeling 157
7.2.5	Evaluation 158
7.2.6	Deployment 158
	2 0 10 10 10 10 10 10 10 10 10 10 10 10 1
	Part III Predicting the Unknown 159
8	Regression 161
8.1	Predictive Performance Estimation 164
8.1.1	Generalization 164
8.1.2	Model Validation 165
8.1.3	Predictive Performance Measures for
0.2.0	Regression 169
8.2	Finding the Parameters of the Model 171
8.2.1	Linear Regression 171
8.2.1.1	Empirical Error 173
8.2.2	The Bias-variance Trade-off 175
8.2.3	Shrinkage Methods 177
8.2.3.1	Ridge Regression 179
8.2.3.2	Lasso Regression 180
8.2.4	Methods that use Linear Combinations of
	Attributes 181
8.2.4.1	Principal Components Regression 181
8.2.4.2	Partial Least Squares Regression 182
8.3	Technique and Model Selection 182
8.4	Final Remarks 183
8.5	Exercises 184
0.0	Exercises 101
9	Classification 187
9.1	Binary Classification 188
9.2	Predictive Performance Measures for Classification 192
9.3	Distance-based Learning Algorithms 199
9.3.1	K-nearest Neighbor Algorithms 199
9.3.2	Case-based Reasoning 202
9.4	Probabilistic Classification Algorithms 203
9.4.1	Logistic Regression Algorithm 205
9.4.2	Naive Bayes Algorithm 207
9.5	Final Remarks 208
9.6	Exercises 210
7.0	Inci cioco 210

10	Additional Predictive Methods 211
10.1	Search-based Algorithms 211
10.1.1	Decision Tree Induction Algorithms 212
10.1.2	Decision Trees for Regression 217
10.1.2.1	Model Trees 218
10.1.2.2	Multivariate Adaptive Regression Splines 219
10.2	Optimization-based Algorithms 221
10.2.1	Artificial Neural Networks 222
10.2.1.1	Backpropagation 224
10.2.1.2	Deep Networks and Deep Learning Algorithms 230
10.2.2	Support Vector Machines 233
10.2.2.1	SVM for Regression 237
10.3	Final Remarks 238
10.4	Exercises 239
11	Advanced Predictive Topics 241
11.1	Ensemble Learning 241
11.1.1	Bagging 243
11.1.2	Random Forests 244
11.1.3	AdaBoost 245
11.2	Algorithm Bias 246
11.3	Non-binary Classification Tasks 248
11.3.1	One-class Classification 248
11.3.2	Multi-class Classification 249
11.3.3	Ranking Classification 250
11.3.4	Multi-label Classification 251
11.3.5	Hierarchical Classification 252
11.4	Advanced Data Preparation Techniques for Prediction 253
11.4.1	Imbalanced Data Classification 253
11.4.2	For Incomplete Target Labeling 254
11.4.2.1	Semi-supervised Learning 254
11.4.2.2	Active Learning 255
11.5	Description and Prediction with Supervised Interpretable
	Techniques 255
11.6	Exercises 256
12	Cheat Sheet and Project on Predictive Analytics 259
12.1	Cheat Sheet on Predictive Analytics 259
12.2	Project on Predictive Analytics 259
12.2.1	Business Understanding 260
12.2.2	Data Understanding 260
12.2.3	Data Preparation 265
12.2.4	Modeling 265
12.2.5	Evaluation 265
12.2.6	Deployment 266

Part IV Popular Data Ana	lytics Applications	267
--------------------------	---------------------	-----

13	Applications for Text, Web and Social Media 269
13.1	Working with Texts 269
13.1.1	Data Acquisition 271
13.1.2	Feature Extraction 271
13.1.2.1	Tokenization 272
13.1.2.2	Stemming 272
	Conversion to Structured Data 275
13.1.2.4	Is the Bag of Words Enough? 276
13.1.3	Remaining Phases 277
13.1.4	Trends 277
13.1.4.1	Sentiment Analysis 278
	Web Mining 278
13.2	Recommender Systems 278
13.2.1	Feedback 279
13.2.2	Recommendation Tasks 280
13.2.3	Recommendation Techniques 281
13.2.3.1	Knowledge-based Techniques 281
13.2.3.2	Content-based Techniques 282
13.2.3.3	Collaborative Filtering Techniques 282
13.2.4	Final Remarks 289
13.3	Social Network Analysis 291
13.3.1	Representing Social Networks 291
13.3.2	Basic Properties of Nodes 294
13.3.2.1	Degree 294
13.3.2.2	Distance 294
13.3.2.3	Closeness 295
13.3.2.4	Betweenness 296
13.3.2.5	Clustering Coefficient 297
13.3.3	Basic and Structural Properties of Networks 297
13.3.3.1	
13.3.3.2	Centralization 297
13.3.3.3	Cliques 299
13.3.3.4	Clustering Coefficient 299
13.3.3.5	Modularity 299
13.3.4	Trends and Final Remarks 299
13.4	Exercises 300

Apendix A: Comprehensive Description of the CRISP-DM Methodology 303 References 311 Index *315*

Preface

We are living in a period of history that will certainly be remembered as one where information began to be instantaneously obtainable, services were tailored to individual criteria, and people did what made them feel good (if it did not put their lives at risk). Every year, machines are able to do more and more things that improve our quality of life. More data is available than ever before, and will become even more so. This is a time when we can extract more information from data than ever before, and benefit more from it.

In different areas of business and in different institutions, new ways to collect data are continuously being created. Old documents are being digitized, new sensors count the number of cars passing along motorways and extract useful information from them, our smartphones are informing us where we are at each moment and what new opportunities are available, and our favorite social networks register to whom we are related or what things we like.

Whatever area we work in, new data is available: data on how students evaluate professors, data on the evolution of diseases and the best treatment options per patient, data on soil, humidity levels and the weather, enabling us to produce more food with better quality, data on the macro economy, our investments and stock market indicators over time, enabling fairer distribution of wealth, data on things we purchase, allowing us to purchase more effectively and at lower cost.

Students in many different domains feel the need to take advantage of the data they have. New courses on data analytics have been proposed in many different programs, from biology to information science, from engineering to economics, from social sciences to agronomy, all over the world.

The first books on data analytics that appeared some years ago were written by data scientists for other data scientists or for data science students. The majority of the people interested in these subjects were computing and statistics students. The books on data analytics were written mainly for them. Nowadays, more and more people are interested in learning data analytics. Students of economics, management, biology, medicine, sociology, engineering, and some other subjects are willing to learn about data analytics. This book intends not only to provide a new, more friendly textbook for computing and statistics students, but also to open data analytics to those students who may know nothing about computing or statistics, but want to learn these subjects in a simple way. Those who have already studied subjects such as statistics will recognize some of the content described in this book, such as descriptive statistics. Students from computing will be familiar with a pseudocode.

After reading this book, it is not expected that you will feel like a data scientist with ability to create new methods, but it is expected that you might feel like a data analytics practitioner, able to drive a data analytics project, using the right methods to solve real problems.

> João Mendes Moreira University of Porto, Porto, Portugal

André C. P. L. F. de Carvalho University of São Paulo, São Carlos, Brazil

Tomáš Horváth Eötvös Loránd University in Budapest Pavol Jozef Šafárik University in Košice October, 2017

Acknowledgments

The authors would like to thank Bruno Almeida Pimentel, Edésio Alcobaça Neto, Everlândio Fernandes, Victor Alexandre Padilha and Victor Hugo Barella for their useful comments.

Over the last several months, we have been in contact with several people from Wiley: Jon Gurstelle, Executive Editor on Statistics; Kathleen Pagliaro, Assistant Editor; Samantha Katherine Clarke and Kshitija Iyer, Project Editors; and Katrina Maceda, Production Editor. To all these wonderful people, we owe a deep sense of gratitude, especially now this project has been completed.

Lastly, we would like to thank our families for their constant love, support, patience, and encouragement.

J. A. T.

Presentational Conventions

Definition The definitions are presented in the format shown here.

Special sections and formats Whenever a method is described, three different sections are presented:

- Assessing and evaluating results: how can we assess the results of a method? How to interpret them? This section is all about answering these questions.
- Setting the hyper-parameters: each method has its own hyper-parameters that must be set. This section explains how to set them.
- Advantages and disadvantages: a table summarizes the positive and negative characteristics of a given method.

About the Companion Website

This book is accompanied by a companion website:

www.wiley.com/go/moreira/dataanalytics

The website includes:

• Presentation slides for instructors

Part I

Introductory Background

1

What Can We Do With Data?

Until recently, researchers working with data analysis were struggling to obtain data for their experiments. Recent advances in the technology of data processing, data storage and data transmission, associated with advanced and intelligent computer software, reducing costs and increasing capacity, have changed this scenario. It is the time of the Internet of Things, where the aim is to have everything or almost everything connected. Data previously produced on paper are now on-line. Each day, a larger quantity of data is generated and consumed. Whenever you place a comment in your social network, upload a photograph, some music or a video, navigate through the Internet, or add a comment to an e-commerce web site, you are contributing to the data increase. Additionally, machines, financial transactions and sensors such as security cameras, are increasingly gathering data from very diverse and widespread sources.

In 2012, it was estimated that, each year, the amount of data available in the world doubles [1]. Another estimate, from 2014, predicted that by 2020 all information will be digitized, eliminated or reinvented in 80% of processes and products of the previous decade [2]. In a third report, from 2015, it was predicted that mobile data traffic will be almost 10 times larger in 2020 [3]. The result of all these rapid increases of data is named by some the "data explosion".

Despite the impression that this can give – that we are drowning in data – there are several benefits from having access to all these data. These data provide a rich source of information that can be transformed into new, useful, valid and human-understandable knowledge. Thus, there is a growing interest in exploring these data to extract this knowledge, using it to support decision making in a wide variety of fields: agriculture, commerce, education, environment, finance, government, industry, medicine, transport and social care. Several companies around the world are realizing the gold mine they have and the potential of these data to support their work, reduce waste and dangerous and tedious work activities, and increase the value of their products and their profits.

The analysis of these data to extract such knowledge is the subject of a vibrant area known as data analytics, or simply "analytics". You can find several definitions of analytics in the literature. The definition adopted here is:

Analytics The science that analyze crude data to extract useful knowledge (patterns) from them.

This process can also include data collection, organization, pre-processing, transformation, modeling and interpretation.

Analytics as a knowledge area involves input from many different areas. The idea of generalizing knowledge from a data sample comes from a branch of statistics known as inductive learning, an area of research with a long history. With the advances of personal computers, the use of computational resources to solve problems of inductive learning become more and more popular. Computational capacity has been used to develop new methods. At the same time, new problems have appeared requiring a good knowledge of computer sciences. For instance, the ability to perform a given task with more computational efficiency has become a subject of study for people working in computational statistics.

In parallel, several researchers have dreamed of being able to reproduce human behavior using computers. These were people from the area of artificial intelligence. They also used statistics for their research but the idea of reproducing human and biological behavior in computers was an important source of motivation. For instance, reproducing how the human brain works with artificial neural networks has been studied since the 1940s; reproducing how ants work with ant colony optimization algorithm since the 1990s. The term machine learning (ML) appeared in this context as the "field of study that gives computers the ability to learn without being explicitly programmed," according to Arthur Samuel in 1959 [4].

In the 1990s, a new term appeared with a different slight meaning: data mining (DM). The 1990s was the decade of the appearance of business intelligence tools as consequence of the data facilities having larger and cheaper capacity. Companies start to collect more and more data, aiming to either solve or improve business operations, for example by detecting frauds with credit cards, by advising the public of road network constraints in cities, or by improving relations with clients using more efficient techniques of relational marketing. The question was of being able to mine the data in order to extract the knowledge necessary for a given task. This is the goal of data mining.

Big Data and Data Science 1.1

In the first years of the 20th century, the term big data has appeared. Big data, a technology for data processing, was initially defined by the "three Vs", although some more Vs have been proposed since. The first three Vs allow us to define a taxonomy of big data. They are: volume, variety and velocity. Volume is concerned with how to store big data: data repositories for large amounts of data. Variety is concerned with how to put together data from different sources. Velocity concerns the ability to deal with data arriving very fast, in streams known as data streams. Analytics is also about discovering knowledge from data streams, going beyond the velocity component of big data.

Another term that has appeared and is sometimes used as a synonym for big data is data science. According to Provost and Fawcett [5], big data are data sets that are too large to be managed by conventional data-processing technologies, requiring the development of new techniques and tools for data storage, processing and transmission. These tools include, for example, MapReduce, Hadoop, Spark and Storm. But data volume is not the only characterization of big data. The word "big" can refer to the number of data sources, to the importance of the data, to the need for new processing techniques, to how fast data arrive, to the combination of different sets of data so they can be analyzed in real time, and its ubiquity, since any company, nonprofit organization or individual has access to data now.

Thus big data is more concerned with technology. It provides a computing environment, not only for analytics, but also for other data processing tasks. These tasks include finance transaction processing, web data processing and georeferenced data processing.

Data science is concerned with the creation of models able to extract patterns from complex data and the use of these models in real-life problems. Data science extracts meaningful and useful knowledge from data, with the support of suitable technologies. It has a close relationship to analytics and data mining. Data science goes beyond data mining by providing a knowledge extraction framework, including statistics and visualization.

Therefore, while big data gives support to data collection and management, data science applies techniques to these data to discover new and useful knowledge: big data collects and data science discovers. Other terms such as knowledge discovery or extraction, pattern recognition, data analysis, data engineering, and several others are also used. The definition we use of data analytics covers all these areas that are used to extract knowledge from data.

1.2 **Big Data Architectures**

As data increase in size, velocity and variety, new computer technologies become necessary. These new technologies, which include hardware and software, must be easily expanded as more data are processed. This property is known as scalability. One way to obtain scalability is by distributing the data processing tasks into several computers, which can be combined into clusters of computers. The reader should not confuse clusters of computers

with clusters produced by clustering techniques, which are techniques from analytics in which a data set is partitioned to find groups within it.

Even if processing power is expanded by combining several computers in a cluster, creating a distributed system, conventional software for distributed systems usually cannot cope with big data. One of the limitations is the efficient distribution of data among the different processing and storage units. To deal with these requirements, new software tools and techniques have been developed.

One of the first techniques developed for big data processing using clusters was MapReduce. MapReduce is a programming model that has two steps: map and reduce. The most famous implementation of MapReduce is called Hadoop.

MapReduce divides the data set into parts – chunks – and stores in the memory of each cluster computer the chunk of the data set needed by this computer to accomplish its processing task. As an example, suppose that you need to calculate the average salary of 1 billion people and you have a cluster with 1000 computers, each with a processing unit and a storage memory. The people can be divided into 1000 chunks – subsets – with data from 1 million people each. Each chunk can be processed independently by one of the computers. The results produced by each these computers (the average salary of 1 million people) can be averaged, returning the final salary average.

To efficiently solve a big data problem, a distributed system must attend the following requirements:

- Make sure that no chunk of data is lost and the whole task is concluded. If
 one or more computers has a failure, their tasks, and the corresponding data
 chunk, must be assumed by another computer in the cluster.
- Repeat the same task, and corresponding data chunk, in more than one cluster computer; this is called redundancy. Thus, if one or more computer fails, the redundant computer carries on with the task.
- Computers that have had faults can return to the cluster again when they are fixed.
- Computers can be easily removed from the cluster or extra ones included in it as the processing demand changes.

A solution incorporating these requirements must hide from the data analyst the details of how the software works, such as how the data chunks and tasks are divided among the cluster computers.

1.3 Small Data

In the opposite direction from big data technologies and methods, there is a movement towards more personal, subjective analysis of chunks of data, termed "small data". Small data is a data set whose volume and format allows its processing and analysis by a person or a small organization. Thus, instead of collecting data from several sources, with different formats, and generated at increasing velocities, creating large data repositories and processing facilities, small data favors the partition of a problem into small packages, which can be analyzed by different people or small groups in a distributed and integrated way.

People are continuously producing small data as they perform their daily activities, be it navigating the web, buying a product in a shop, undergoing medical examinations and using apps in their mobiles. When these data are collected to be stored and processed in large data servers they become big data. To be characterized as small data, a data set must have a size that allows its full understanding by an user.

The type of knowledge sought in big and small data is also different, with the first looking for correlations and the second for causality relations. While big data provide tools that allow companies to understand their customers, small data tools try to help customers to understand themselves. Thus, big data is concerned with customers, products and services, and small data is concerned with the individuals that produced the data.

1.4 What is Data?

But what is data about? Data, in the information age, are a large set of bits encoding numbers, texts, images, sounds, videos, and so on. Unless we add information to data, they are meaningless. When we add information, giving a meaning to them, these data become knowledge. But before data become knowledge, typically, they pass through several steps where they are still referred to as data, despite being a bit more organized; that is, they have some information associated with them.

Let us see the example of data collected from a private list of acquaintances or contacts.

Information as presented in Table 1.1, usually referred to as tabular data, is characterized by the way data are organized. In tabular data, data are organized in rows and columns, where each column represents a characteristic of the data and each row represents an occurrence of the data. A column is referred to as an attribute or, with the same meaning, a feature, while a row is referred to as an instance, or with the same meaning, an object.

Instance or Object Examples of the concept we want to characterize.

Example 1.1 In the example in Table 1.1, we intend to characterize people in our private contact list. Each member is, in this case, an instance or object. It corresponds to a row of the table.

Attribute or Feature Attributes, also called features, are characteristics of the instances.

Tabl	e 1.1	Data set o	f our	private	contact l	ist.
------	-------	------------	-------	---------	-----------	------

Contact	Age	Educational level	Company
Andrew	55	1.0	Good
Bernhard	43	2.0	Good
Carolina	37	5.0	Bad
Dennis	82	3.0	Good
Eve	23	3.2	Bad
Fred	46	5.0	Good
Gwyneth	38	4.2	Bad
Hayden	50	4.0	Bad
Irene	29	4.5	Bad
James	42	4.1	Good
Kevin	35	4.5	Bad
Lea	38	2.5	Good
Marcus	31	4.8	Bad
Nigel	71	2.3	Good

Example 1.2 In Table 1.1, contact, age, education level and company are four different attributes.

The majority of the chapters in this book expect the data to be in tabular format; that is, already organized by rows and columns, each row representing an instance and each column representing an attribute. However, a table can be organized differently, having the instances per column and the attributes per row.

There are, however, data that are not possible to represent in a single table.

Example 1.3 As an example, if some of the contacts are relatives of other contacts, a second table, as shown in Table 1.2, representing the family relationships, would be necessary. You should note that each person referred to in Table 1.2 also exists in Table 1.1, i.e., there are relations between attributes of different tables.

Data sets represented by several tables, making clear the relations between these tables, are called relational data sets. This information is easily handled using relational databases. In this book, only simple forms of relational data will be used. This is discussed in each chapter whenever necessary.