



Python Data Analytics

With Pandas, NumPy, and Matplotlib

—
Second Edition

—
Fabio Nelli

Apress®

Python Data Analytics

With Pandas, NumPy,
and Matplotlib

Second Edition

Fabio Nelli

Apress®

Python Data Analytics

Fabio Nelli
Rome, Italy

ISBN-13 (pbk): 978-1-4842-3912-4
<https://doi.org/10.1007/978-1-4842-3913-1>

ISBN-13 (electronic): 978-1-4842-3913-1

Library of Congress Control Number: 2018957991

Copyright © 2018 by Fabio Nelli

This work is subject to copyright. All rights are reserved by the Publisher, whether the whole or part of the material is concerned, specifically the rights of translation, reprinting, reuse of illustrations, recitation, broadcasting, reproduction on microfilms or in any other physical way, and transmission or information storage and retrieval, electronic adaptation, computer software, or by similar or dissimilar methodology now known or hereafter developed.

Trademarked names, logos, and images may appear in this book. Rather than use a trademark symbol with every occurrence of a trademarked name, logo, or image we use the names, logos, and images only in an editorial fashion and to the benefit of the trademark owner, with no intention of infringement of the trademark.

The use in this publication of trade names, trademarks, service marks, and similar terms, even if they are not identified as such, is not to be taken as an expression of opinion as to whether or not they are subject to proprietary rights.

While the advice and information in this book are believed to be true and accurate at the date of publication, neither the authors nor the editors nor the publisher can accept any legal responsibility for any errors or omissions that may be made. The publisher makes no warranty, express or implied, with respect to the material contained herein.

Managing Director, Apress Media LLC: Welmoed Spahr
Acquisitions Editor: Todd Green
Development Editor: James Markham
Coordinating Editor: Jill Balzano

Cover image designed by Freepik (www.freepik.com)

Distributed to the book trade worldwide by Springer Science+Business Media New York, 233 Spring Street, 6th Floor, New York, NY 10013. Phone 1-800-SPRINGER, fax (201) 348-4505, e-mail orders-ny@springer-sbm.com, or visit www.springeronline.com. Apress Media, LLC is a California LLC and the sole member (owner) is Springer Science + Business Media Finance Inc (SSBM Finance Inc). SSBM Finance Inc is a **Delaware** corporation.

For information on translations, please e-mail rights@apress.com, or visit <http://www.apress.com/rights-permissions>.

Apress titles may be purchased in bulk for academic, corporate, or promotional use. eBook versions and licenses are also available for most titles. For more information, reference our Print and eBook Bulk Sales web page at <http://www.apress.com/bulk-sales>.

Any source code or other supplementary material referenced by the author in this book is available to readers on GitHub via the book's product page, located at www.apress.com/9781484239124. For more detailed information, please visit <http://www.apress.com/source-code>.

Printed on acid-free paper

*“Science leads us forward in knowledge, but only analysis
makes us more aware”*

*This book is dedicated to all those who are constantly
looking for awareness*

Table of Contents

About the Author	xvii
About the Technical Reviewer	xix
Chapter 1: An Introduction to Data Analysis	1
Data Analysis	1
Knowledge Domains of the Data Analyst	3
Computer Science	3
Mathematics and Statistics	4
Machine Learning and Artificial Intelligence	5
Professional Fields of Application.....	5
Understanding the Nature of the Data	5
When the Data Become Information.....	6
When the Information Becomes Knowledge	6
Types of Data.....	6
The Data Analysis Process	6
Problem Definition.....	8
Data Extraction	9
Data Preparation.....	10
Data Exploration/Visualization.....	10
Predictive Modeling.....	12
Model Validation	13
Deployment	13
Quantitative and Qualitative Data Analysis	14
Open Data	15
Python and Data Analysis.....	17
Conclusions.....	17

TABLE OF CONTENTS

- Chapter 2: Introduction to the Python World 19**
 - Python—The Programming Language..... 19
 - Python—The Interpreter..... 21
 - Python 2 and Python 3 23
 - Installing Python 23
 - Python Distributions 24
 - Using Python..... 26
 - Writing Python Code 28
 - IPython..... 35
 - PyPI—The Python Package Index..... 39
 - The IDEs for Python 40
 - SciPy 46
 - NumPy 47
 - Pandas..... 47
 - matplotlib 48
 - Conclusions..... 48
- Chapter 3: The NumPy Library 49**
 - NumPy: A Little History..... 49
 - The NumPy Installation 50
 - Ndarray: The Heart of the Library..... 50
 - Create an Array 52
 - Types of Data 53
 - The dtype Option 54
 - Intrinsic Creation of an Array 55
 - Basic Operations 57
 - Arithmetic Operators 57
 - The Matrix Product 59
 - Increment and Decrement Operators 60
 - Universal Functions (ufunc) 61
 - Aggregate Functions 62

Indexing, Slicing, and Iterating	62
Indexing	63
Slicing	65
Iterating an Array	67
Conditions and Boolean Arrays	69
Shape Manipulation	70
Array Manipulation	71
Joining Arrays	71
Splitting Arrays	72
General Concepts	74
Copies or Views of Objects	75
Vectorization	76
Broadcasting	76
Structured Arrays	79
Reading and Writing Array Data on Files	82
Loading and Saving Data in Binary Files	82
Reading Files with Tabular Data	83
Conclusions	84
Chapter 4: The pandas Library—An Introduction	87
pandas: The Python Data Analysis Library	87
Installation of pandas	88
Installation from Anaconda	88
Installation from PyPI	89
Installation on Linux	90
Installation from Source	90
A Module Repository for Windows	90
Testing Your pandas Installation	91
Getting Started with pandas	92
Introduction to pandas Data Structures	92
The Series	93

TABLE OF CONTENTS

- The DataFrame 102
- The Index Objects 112
- Other Functionalities on Indexes 114
 - Reindexing 114
 - Dropping 117
 - Arithmetic and Data Alignment 118
- Operations Between Data Structures 120
 - Flexible Arithmetic Methods 120
 - Operations Between DataFrame and Series 121
- Function Application and Mapping 122
 - Functions by Element 123
 - Functions by Row or Column 123
 - Statistics Functions 125
- Sorting and Ranking 126
- Correlation and Covariance 129
- “Not a Number” Data 131
 - Assigning a NaN Value 131
 - Filtering Out NaN Values 132
 - Filling in NaN Occurrences 133
- Hierarchical Indexing and Leveling 134
 - Reordering and Sorting Levels 137
 - Summary Statistic by Level 138
- Conclusions 139
- Chapter 5: pandas: Reading and Writing Data 141**
 - I/O API Tools 141
 - CSV and Textual Files 142
 - Reading Data in CSV or Text Files 143
 - Using RegExp to Parse TXT Files 146
 - Reading TXT Files Into Parts 148
 - Writing Data in CSV 150

Reading and Writing HTML Files	152
Writing Data in HTML.....	153
Reading Data from an HTML File.....	155
Reading Data from XML	157
Reading and Writing Data on Microsoft Excel Files	159
JSON Data.....	162
The Format HDF5	166
Pickle—Python Object Serialization	168
Serialize a Python Object with cPickle	168
Pickling with pandas	169
Interacting with Databases	170
Loading and Writing Data with SQLite3	171
Loading and Writing Data with PostgreSQL.....	174
Reading and Writing Data with a NoSQL Database: MongoDB.....	178
Conclusions.....	180
Chapter 6: pandas in Depth: Data Manipulation	181
Data Preparation	181
Merging	182
Concatenating.....	188
Combining	191
Pivoting.....	193
Removing.....	196
Data Transformation.....	197
Removing Duplicates.....	198
Mapping.....	199
Discretization and Binning	204
Detecting and Filtering Outliers.....	209
Permutation	210
Random Sampling	211

TABLE OF CONTENTS

- String Manipulation..... 212
 - Built-in Methods for String Manipulation 212
 - Regular Expressions 214
- Data Aggregation 217
 - GroupBy 218
 - A Practical Example..... 219
 - Hierarchical Grouping 220
- Group Iteration 222
 - Chain of Transformations..... 222
 - Functions on Groups..... 224
- Advanced Data Aggregation..... 225
- Conclusions..... 229
- Chapter 7: Data Visualization with matplotlib 231**
 - The matplotlib Library 231
 - Installation 233
 - The IPython and IPython QtConsole 233
 - The matplotlib Architecture..... 235
 - Backend Layer 236
 - Artist Layer 236
 - Scripting Layer (pyplot) 238
 - pylab and pyplot 238
 - pyplot..... 239
 - A Simple Interactive Chart..... 239
 - The Plotting Window 241
 - Set the Properties of the Plot 243
 - matplotlib and NumPy 246
 - Using the kwargs 248
 - Working with Multiple Figures and Axes 249
 - Adding Elements to the Chart 251
 - Adding Text 251

Adding a Grid	256
Adding a Legend.....	257
Saving Your Charts	260
Saving the Code.....	260
Converting Your Session to an HTML File	262
Saving Your Chart Directly as an Image.....	264
Handling Date Values	264
Chart Typology.....	267
Line Charts	267
Line Charts with pandas.....	276
Histograms.....	277
Bar Charts	278
Horizontal Bar Charts.....	281
Multiserial Bar Charts.....	282
Multiseries Bar Charts with pandas Dataframe.....	285
Multiseries Stacked Bar Charts	286
Stacked Bar Charts with a pandas Dataframe	290
Other Bar Chart Representations.....	291
Pie Charts.....	292
Pie Charts with a pandas Dataframe	296
Advanced Charts	297
Contour Plots	297
Polar Charts	299
The mplot3d Toolkit.....	302
3D Surfaces	302
Scatter Plots in 3D.....	304
Bar Charts in 3D	306
Multi-Panel Plots.....	307
Display Subplots Within Other Subplots	307
Grids of Subplots	309
Conclusions.....	312

TABLE OF CONTENTS

- Chapter 8: Machine Learning with scikit-learn 313**
 - The scikit-learn Library..... 313
 - Machine Learning 313
 - Supervised and Unsupervised Learning..... 314
 - Training Set and Testing Set..... 315
 - Supervised Learning with scikit-learn..... 315
 - The Iris Flower Dataset..... 316
 - The PCA Decomposition 320
 - K-Nearest Neighbors Classifier..... 322
 - Diabetes Dataset..... 327
 - Linear Regression: The Least Square Regression..... 328
 - Support Vector Machines (SVMs)..... 334
 - Support Vector Classification (SVC)..... 334
 - Nonlinear SVC..... 339
 - Plotting Different SVM Classifiers Using the Iris Dataset 342
 - Support Vector Regression (SVR)..... 345
 - Conclusions..... 347

- Chapter 9: Deep Learning with TensorFlow 349**
 - Artificial Intelligence, Machine Learning, and Deep Learning 349
 - Artificial intelligence..... 350
 - Machine Learning Is a Branch of Artificial Intelligence 351
 - Deep Learning Is a Branch of Machine Learning..... 351
 - The Relationship Between Artificial Intelligence, Machine Learning, and Deep Learning... 351
 - Deep Learning..... 352
 - Neural Networks and GPUs 352
 - Data Availability: Open Data Source, Internet of Things, and Big Data 353
 - Python 354
 - Deep Learning Python Frameworks 354
 - Artificial Neural Networks..... 355
 - How Artificial Neural Networks Are Structured 355
 - Single Layer Perceptron (SLP)..... 357

Multi Layer Perceptron (MLP)	360
Correspondence Between Artificial and Biological Neural Networks	361
TensorFlow	362
TensorFlow: Google's Framework	362
TensorFlow: Data Flow Graph	362
Start Programming with TensorFlow	363
Installing TensorFlow	363
Programming with the IPython QtConsole	364
The Model and Sessions in TensorFlow	364
Tensors	366
Operation on Tensors	370
Single Layer Perceptron with TensorFlow	371
Before Starting	372
Data To Be Analyzed	372
The SLP Model Definition	374
Learning Phase	378
Test Phase and Accuracy Calculation	383
Multi Layer Perceptron (with One Hidden Layer) with TensorFlow	386
The MLP Model Definition	387
Learning Phase	389
Test Phase and Accuracy Calculation	395
Multi Layer Perceptron (with Two Hidden Layers) with TensorFlow	397
Test Phase and Accuracy Calculation	402
Evaluation of Experimental Data	404
Conclusions	407
Chapter 10: An Example— Meteorological Data	409
A Hypothesis to Be Tested: The Influence of the Proximity of the Sea	409
The System in the Study: The Adriatic Sea and the Po Valley	410
Finding the Data Source	414
Data Analysis on Jupyter Notebook	415
Analysis of Processed Meteorological Data	421

TABLE OF CONTENTS

- The RoseWind 436
 - Calculating the Mean Distribution of the Wind Speed 441
- Conclusions..... 443
- Chapter 11: Embedding the JavaScript D3 Library in the IPython Notebook 445**
- The Open Data Source for Demographics 445
- The JavaScript D3 Library 449
- Drawing a Clustered Bar Chart 454
- The Choropleth Maps 459
- The Choropleth Map of the U.S. Population in 2014..... 464
- Conclusions..... 471
- Chapter 12: Recognizing Handwritten Digits..... 473**
- Handwriting Recognition..... 473
- Recognizing Handwritten Digits with scikit-learn..... 474
- The Digits Dataset..... 475
- Learning and Predicting..... 478
- Recognizing Handwritten Digits with TensorFlow..... 480
- Learning and Predicting..... 482
- Conclusions..... 486
- Chapter 13: Textual Data Analysis with NLTK 487**
- Text Analysis Techniques 487
 - The Natural Language Toolkit (NLTK) 488
 - Import the NLTK Library and the NLTK Downloader Tool 489
 - Search for a Word with NLTK..... 493
 - Analyze the Frequency of Words 494
 - Selection of Words from Text..... 497
 - Bigrams and Collocations..... 498
- Use Text on the Network 500
 - Extract the Text from the HTML Pages 501
 - Sentimental Analysis 502
- Conclusions..... 506

Chapter 14: Image Analysis and Computer Vision with OpenCV	507
Image Analysis and Computer Vision	507
OpenCV and Python.....	508
OpenCV and Deep Learning	509
Installing OpenCV.....	509
First Approaches to Image Processing and Analysis.....	509
Before Starting	510
Load and Display an Image	510
Working with Images.....	512
Save the New Image.....	514
Elementary Operations on Images	514
Image Blending.....	520
Image Analysis	521
Edge Detection and Image Gradient Analysis	522
Edge Detection	522
The Image Gradient Theory	523
A Practical Example of Edge Detection with the Image Gradient Analysis.....	525
A Deep Learning Example: The Face Detection.....	532
Conclusions.....	535
Appendix A: Writing Mathematical Expressions with LaTeX	537
With matplotlib.....	537
With IPython Notebook in a Markdown Cell.....	537
With IPython Notebook in a Python 2 Cell.....	538
Subscripts and Superscripts.....	538
Fractions, Binomials, and Stacked Numbers	538
Radicals	539
Fonts	539
Accents	540

TABLE OF CONTENTS

Appendix B: Open Data Sources 549

- Political and Government Data..... 549
- Health Data 550
- Social Data 550
- Miscellaneous and Public Data Sets 551
- Financial Data 552
- Climatic Data..... 552
- Sports Data 553
- Publications, Newspapers, and Books 553
- Musical Data 553

Index..... 555

About the Author

Fabio Nelli is a data scientist and Python consultant, designing and developing Python applications for data analysis and visualization. He has experience with the scientific world, having performed various data analysis roles in pharmaceutical chemistry for private research companies and universities. He has been a computer consultant for many years at IBM, EDS, and Hewlett-Packard, along with several banks and insurance companies. He has an organic chemistry master's degree and a bachelor's degree in information technologies and automation systems, with many years of experience in life sciences (as as Tech Specialist at Beckman Coulter, Tecan, Sciex).

For further info and other examples, visit his page at <https://www.meccanismocomplesso.org> and the GitHub page <https://github.com/meccanismocomplesso>.

About the Technical Reviewer



Raul Samayoa is a senior software developer and machine learning specialist with many years of experience in the financial industry. An MSc graduate from the Georgia Institute of Technology, he's never met a neural network or dataset he did not like. He's fond of evangelizing the use of DevOps tools for data science and software development.

Raul enjoys the energy of his hometown of Toronto, Canada, where he runs marathons, volunteers as a technology instructor with the University of Toronto coders, and likes to work with data in Python and R.

CHAPTER 1

An Introduction to Data Analysis

In this chapter, you begin to take the first steps in the world of data analysis, learning in detail about all the concepts and processes that make up this discipline. The concepts discussed in this chapter are helpful background for the following chapters, where these concepts and procedures will be applied in the form of Python code, through the use of several libraries that will be discussed in just as many chapters.

Data Analysis

In a world increasingly centralized around information technology, huge amounts of data are produced and stored each day. Often these data come from automatic detection systems, sensors, and scientific instrumentation, or you produce them daily and unconsciously every time you make a withdrawal from the bank or make a purchase, when you record various blogs, or even when you post on social networks.

But what are the data? The data actually are not information, at least in terms of their form. In the formless stream of bytes, at first glance it is difficult to understand their essence if not strictly the number, word, or time that they report. Information is actually the result of processing, which, taking into account a certain dataset, extracts some conclusions that can be used in various ways. This process of extracting information from raw data is called *data analysis*.

The purpose of data analysis is to extract information that is not easily deducible but that, when understood, leads to the possibility of carrying out studies on the mechanisms of the systems that have produced them, thus allowing you to forecast possible responses of these systems and their evolution in time.

Starting from a simple methodical approach on data protection, data analysis has become a real discipline, leading to the development of real methodologies generating *models*. The model is in fact the translation into a mathematical form of a system placed under study. Once there is a mathematical or logical form that can describe system responses under different levels of precision, you can then make predictions about its development or response to certain inputs. Thus the aim of data analysis is not the model, but the quality of its *predictive power*.

The predictive power of a model depends not only on the quality of the modeling techniques but also on the ability to choose a good dataset upon which to build the entire data analysis process. So the *search for data*, their *extraction*, and their subsequent *preparation*, while representing preliminary activities of an analysis, also belong to data analysis itself, because of their importance in the success of the results.

So far we have spoken of data, their handling, and their processing through calculation procedures. In parallel to all stages of processing of data analysis, various methods of *data visualization* have been developed. In fact, to understand the data, both individually and in terms of the role they play in the entire dataset, there is no better system than to develop the techniques of graphic representation capable of transforming information, sometimes implicitly hidden, in figures, which help you more easily understand their meaning. Over the years lots of display modes have been developed for different modes of data display: the *charts*.

At the end of the data analysis process, you will have a model and a set of graphical displays and then you will be able to predict the responses of the system under study; after that, you will move to the test phase. The model will be tested using another set of data for which you know the system response. These data are, however, not used to define the predictive model. Depending on the ability of the model to replicate real observed responses, you will have an error calculation and knowledge of the validity of the model and its operating limits.

These results can be compared with any other models to understand if the newly created one is more efficient than the existing ones. Once you have assessed that, you can move to the last phase of data analysis—*deployment*. This consists of implementing the results produced by the analysis, namely, implementing the decisions to be taken based on the predictions generated by the model and the associated risks.

Data analysis is well suited to many professional activities. So, knowledge of it and how it can be put into practice is relevant. It allows you to test hypotheses and to understand more deeply the systems analyzed.

Knowledge Domains of the Data Analyst

Data analysis is basically a discipline suitable to the study of problems that may occur in several fields of applications. Moreover, data analysis includes many tools and methodologies that require good knowledge of computing, mathematical, and statistical concepts.

A good data analyst must be able to move and act in many different disciplinary areas. Many of these disciplines are the basis of the methods of data analysis, and proficiency in them is almost necessary. Knowledge of other disciplines is necessary depending on the area of application and study of the particular data analysis project you are about to undertake, and, more generally, sufficient experience in these areas can help you better understand the issues and the type of data needed.

Often, regarding major problems of data analysis, it is necessary to have an interdisciplinary team of experts who can contribute in the best possible way in their respective fields of competence. Regarding smaller problems, a good analyst must be able to recognize problems that arise during data analysis, inquire to determine which disciplines and skills are necessary to solve these problems, study these disciplines, and maybe even ask the most knowledgeable people in the sector. In short, the analyst must be able to know how to search not only for data, but also for information on how to treat that data.

Computer Science

Knowledge of computer science is a basic requirement for any data analyst. In fact, only when you have good knowledge of and experience in computer science can you efficiently manage the necessary tools for data analysis. In fact, every step concerning data analysis involves using calculation software (such as IDL, MATLAB, etc.) and programming languages (such as C ++, Java, and Python).

The large amount of data available today, thanks to information technology, requires specific skills in order to be managed as efficiently as possible. Indeed, data research and extraction require knowledge of these various formats. The data are structured and stored in files or database tables with particular formats. XML, JSON, or simply XLS or CSV files, are now the common formats for storing and collecting data, and many applications allow you to read and manage the data stored on them. When it comes to extracting data contained in a database, things are not so immediate, but you need to know the SQL query language or use software specially developed for the extraction of data from a given database.

Moreover, for some specific types of data research, the data are not available in an explicit format, but are present in text files (documents and log files) or web pages, and shown as charts, measures, number of visitors, or HTML tables. This requires specific technical expertise for the parsing and the eventual extraction of these data (called *web scraping*).

So, knowledge of information technology is necessary to know how to use the various tools made available by contemporary computer science, such as applications and programming languages. These tools, in turn, are needed to perform data analysis and data visualization.

The purpose of this book is to provide all the necessary knowledge, as far as possible, regarding the development of methodologies for data analysis. The book uses the Python programming language and specialized libraries that provide a decisive contribution to the performance of all the steps constituting data analysis, from data research to data mining, to publishing the results of the predictive model.

Mathematics and Statistics

As you will see throughout the book, data analysis requires a lot of complex math during the treatment and processing of data. You need to be competent in all of this, at least to understand what you are doing. Some familiarity with the main statistical concepts is also necessary because all the methods that are applied in the analysis and interpretation of data are based on these concepts. Just as you can say that computer science gives you the tools for data analysis, so you can say that the statistics provide the concepts that form the basis of data analysis.

This discipline provides many tools to the analyst, and a good knowledge of how to best use them requires years of experience. Among the most commonly used statistical techniques in data analysis are

- Bayesian methods
- Regression
- Clustering

Having to deal with these cases, you'll discover how mathematics and statistics are closely related. Thanks to the special Python libraries covered in this book, you will be able to manage and handle them.

Machine Learning and Artificial Intelligence

One of the most advanced tools that falls in the data analysis camp is machine learning. In fact, despite the data visualization and techniques such as clustering and regression, which should help you find information about the dataset, during this phase of research, you may often prefer to use special procedures that are highly specialized in searching patterns within the dataset.

Machine learning is a discipline that uses a whole series of procedures and algorithms that analyze the data in order to recognize patterns, clusters, or trends and then extracts useful information for data analysis in an automated way.

This discipline is increasingly becoming a fundamental tool of data analysis, and thus knowledge of it, at least in general, is of fundamental importance to the data analyst.

Professional Fields of Application

Another very important point is the domain of competence of the data (its source—biology, physics, finance, materials testing, statistics on population, etc.). In fact, although analysts have had specialized preparation in the field of statistics, they must also be able to document the source of the data, with the aim of perceiving and better understanding the mechanisms that generated the data. In fact, the data are not simple strings or numbers; they are the expression, or rather the measure, of any parameter observed. Thus, better understanding where the data came from can improve their interpretation. Often, however, this is too costly for data analysts, even ones with the best intentions, and so it is good practice to find consultants or key figures to whom you can pose the right questions.

Understanding the Nature of the Data

The object of study of data analysis is basically the data. The data then will be the key player in all processes of data analysis. The data constitute the raw material to be processed, and thanks to their processing and analysis, it is possible to extract a variety of information in order to increase the level of knowledge of the system under study, that is, one from which the data came.

When the Data Become Information

Data are the events recorded in the world. Anything that can be measured or categorized can be converted into data. Once collected, these data can be studied and analyzed, both to understand the nature of the events and very often also to make predictions or at least to make informed decisions.

When the Information Becomes Knowledge

You can speak of knowledge when the information is converted into a set of rules that helps you better understand certain mechanisms and therefore make predictions on the evolution of some events.

Types of Data

Data can be divided into two distinct categories:

- Categorical (nominal and ordinal)
- Numerical (discrete and continuous)

Categorical data are values or observations that can be divided into groups or categories. There are two types of categorical values: *nominal* and *ordinal*. A nominal variable has no intrinsic order that is identified in its category. An ordinal variable instead has a predetermined order.

Numerical data are values or observations that come from measurements. There are two types of numerical values: *discrete* and *continuous* numbers. Discrete values can be counted and are distinct and separated from each other. Continuous values, on the other hand, are values produced by measurements or observations that assume any value within a defined range.

The Data Analysis Process

Data analysis can be described as a process consisting of several steps in which the raw data are transformed and processed in order to produce data visualizations and make predictions thanks to a mathematical model based on the collected data. Then, data

analysis is nothing more than a sequence of steps, each of which plays a key role in the subsequent ones. So, data analysis is schematized as a process chain consisting of the following sequence of stages:

- Problem definition
- Data extraction
- Data preparation - Data cleaning
- Data preparation - Data transformation
- Data exploration and visualization
- Predictive modeling
- Model validation/test
- Deploy - Visualization and interpretation of results
- Deploy - Deployment of the solution

Figure 1-1 shows a schematic representation of all the processes involved in the data analysis.

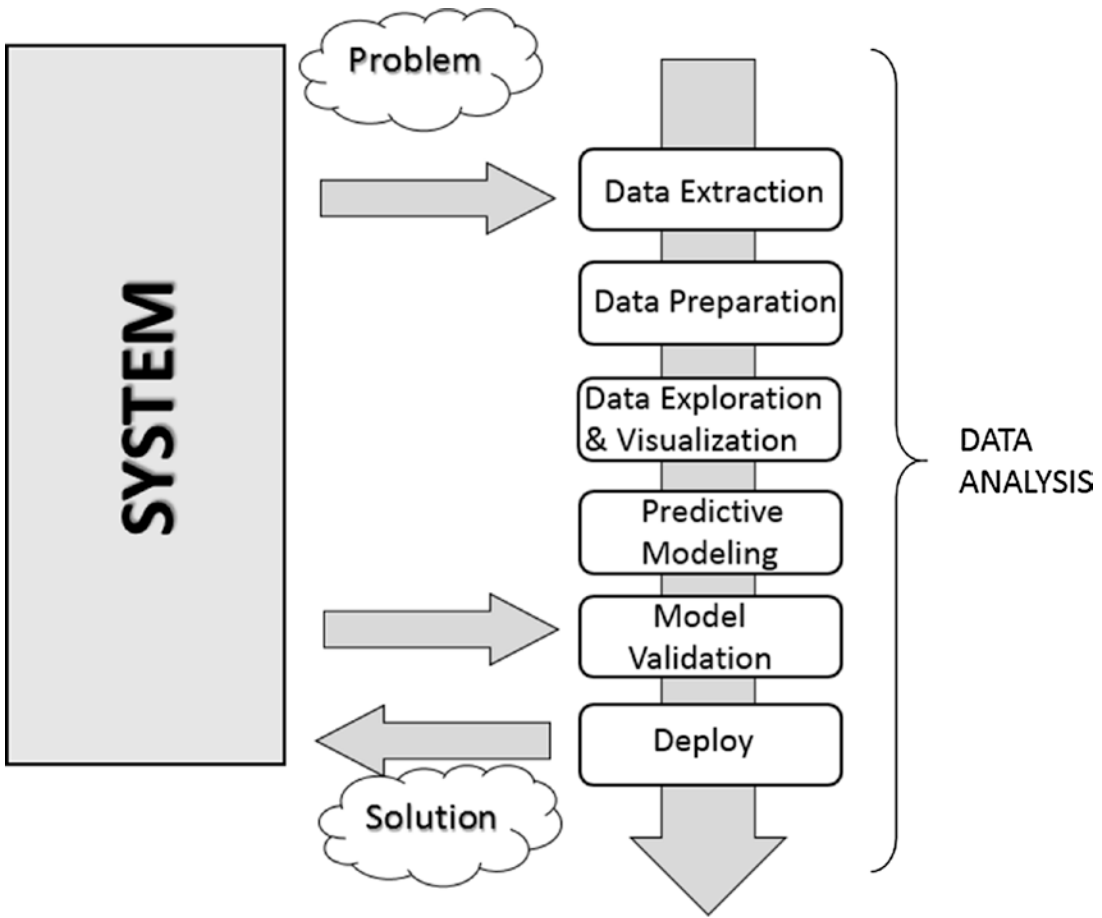


Figure 1-1. The data analysis process

Problem Definition

The process of data analysis actually begins long before the collection of raw data. In fact, data analysis always starts with a problem to be solved, which needs to be defined.

The problem is defined only after you have focused the system you want to study; this may be a mechanism, an application, or a process in general. Generally this study can be in order to better understand its operation, but in particular the study will be designed to understand the principles of its behavior in order to be able to make predictions or choices (defined as an informed choice).

The definition step and the corresponding documentation (*deliverables*) of the scientific problem or business are both very important in order to focus the entire analysis strictly on getting results. In fact, a comprehensive or exhaustive study of the

system is sometimes complex and you do not always have enough information to start with. So the definition of the problem and especially its planning can determine the guidelines to follow for the whole project.

Once the problem has been defined and documented, you can move to the *project planning* stage of data analysis. Planning is needed to understand which professionals and resources are necessary to meet the requirements to carry out the project as efficiently as possible. So you're going to consider the issues in the area involving the resolution of the problem. You will look for specialists in various areas of interest and install the software needed to perform data analysis.

Also during the planning phase, you choose an effective team. Generally, these teams should be cross-disciplinary in order to solve the problem by looking at the data from different perspectives. So, building a good team is certainly one of the key factors leading to success in data analysis.

Data Extraction

Once the problem has been defined, the first step is to obtain the data in order to perform the analysis. The data must be chosen with the basic purpose of building the predictive model, and so data selection is crucial for the success of the analysis as well. The sample data collected must reflect as much as possible the real world, that is, how the system responds to stimuli from the real world. For example, if you're using huge datasets of raw data and they are not collected competently, these may portray false or unbalanced situations.

Thus, poor choice of data, or even performing analysis on a dataset that's not perfectly representative of the system, will lead to models that will move away from the system under study.

The search and retrieval of data often require a form of intuition that goes beyond mere technical research and data extraction. This process also requires a careful understanding of the nature and form of the data, which only good experience and knowledge in the problem's application field can provide.

Regardless of the quality and quantity of data needed, another issue is using the best *data sources*.

If the studio environment is a laboratory (technical or scientific) and the data generated are experimental, then in this case the data source is easily identifiable. In this case, the problems will be only concerning the experimental setup.

But it is not possible for data analysis to reproduce systems in which data are gathered in a strictly experimental way in every field of application. Many fields require searching for data from the surrounding world, often relying on external experimental data, or even more often collecting them through interviews or surveys. So in these cases, finding a good data source that is able to provide all the information you need for data analysis can be quite challenging. Often it is necessary to retrieve data from multiple data sources to supplement any shortcomings, to identify any discrepancies, and to make the dataset as general as possible.

When you want to get the data, a good place to start is the Web. But most of the data on the Web can be difficult to capture; in fact, not all data are available in a file or database, but might be content that is inside HTML pages in many different formats. To this end, a methodology called *web scraping* allows the collection of data through the recognition of specific occurrence of HTML tags within web pages. There is software specifically designed for this purpose, and once an occurrence is found, it extracts the desired data. Once the search is complete, you will get a list of data ready to be subjected to data analysis.

Data Preparation

Among all the steps involved in data analysis, data preparation, although seemingly less problematic, in fact requires more resources and more time to be completed. Data are often collected from different data sources, each of which will have data in it with a different representation and format. So, all of these data will have to be prepared for the process of data analysis.

The preparation of the data is concerned with obtaining, cleaning, normalizing, and transforming data into an optimized dataset, that is, in a prepared format that's normally tabular and is suitable for the methods of analysis that have been scheduled during the design phase.

Many potential problems can arise, including invalid, ambiguous, or missing values, replicated fields, and out-of-range data.

Data Exploration/Visualization

Exploring the data involves essentially searching the data in a graphical or statistical presentation in order to find patterns, connections, and relationships. Data visualization is the best tool to highlight possible patterns.

In recent years, data visualization has been developed to such an extent that it has become a real discipline in itself. In fact, numerous technologies are utilized exclusively to display data, and many display types are applied to extract the best possible information from a dataset.

Data exploration consists of a preliminary examination of the data, which is important for understanding the type of information that has been collected and what it means. In combination with the information acquired during the definition problem, this categorization will determine which method of data analysis will be most suitable for arriving at a model definition.

Generally, this phase, in addition to a detailed study of charts through the visualization data, may consist of one or more of the following activities:

- Summarizing data
- Grouping data
- Exploring the relationship between the various attributes
- Identifying patterns and trends
- Constructing regression models
- Constructing classification models

Generally, data analysis requires summarizing statements regarding the data to be studied. *Summarization* is a process by which data are reduced to interpretation without sacrificing important information.

Clustering is a method of data analysis that is used to find groups united by common attributes (also called *grouping*).

Another important step of the analysis focuses on the *identification* of relationships, trends, and anomalies in the data. In order to find this kind of information, you often have to resort to the tools as well as perform another round of data analysis, this time on the data visualization itself.

Other methods of data mining, such as decision trees and association rules, automatically extract important facts or rules from the data. These approaches can be used in parallel with data visualization to uncover relationships between the data.

Predictive Modeling

Predictive modeling is a process used in data analysis to create or choose a suitable statistical model to predict the probability of a result.

After exploring the data, you have all the information needed to develop the mathematical model that encodes the relationship between the data. These models are useful for understanding the system under study, and in a specific way they are used for two main purposes. The first is to make predictions about the data values produced by the system; in this case, you will be dealing with *regression models*. The second purpose is to classify new data products, and in this case, you will be using *classification models* or *clustering models*. In fact, it is possible to divide the models according to the type of result they produce:

- *Classification models*: If the result obtained by the model type is categorical.
- *Regression models*: If the result obtained by the model type is numeric.
- *Clustering models*: If the result obtained by the model type is descriptive.

Simple methods to generate these models include techniques such as linear regression, logistic regression, classification and regression trees, and k-nearest neighbors. But the methods of analysis are numerous, and each has specific characteristics that make it excellent for some types of data and analysis. Each of these methods will produce a specific model, and then their choice is relevant to the nature of the product model.

Some of these models will provide values corresponding to the real system and according to their structure. They will explain some characteristics of the system under study in a simple and clear way. Other models will continue to give good predictions, but their structure will be no more than a “black box” with limited ability to explain characteristics of the system.