

SEBASTIAN FITZEK

**Introduction to
DATA SCIENCE**

**A PYTHON Path for a
Non-computer Scientist**

Redactor: ANCA MILU-VAIDSEGAN
Tehnoredactor: CRISTIAN LUPEANU

Toate drepturile asupra prezentei ediții aparțin
Editurii COMUNICARE.RO, 2023.

Editura COMUNICARE.RO este departament în cadrul
Școlii Naționale de Studii Politice și Administrative,
Facultatea de Comunicare și Relații Publice.

Editura COMUNICARE.RO

SNSPA, Facultatea de Comunicare și Relații Publice
Str. Povernei, nr. 6, sector 1
010643, București
România
Tel.: 0372.177.150
www.edituracomunicare.ro
e-mail: editura@comunicare.ro

ISBN 978-973-711-642-0 (ediție electronică)

CONTENTS

<i>Biographical data</i>	/ 11
<i>Acknowledgments</i>	/ 12
<i>Foreword</i>	/ 13
<i>An introduction to Data Science from the perspective of the communication and public relations specialist</i>	/ 17

Part I: PROGRAMMING / 23

Python fundamentals in the introduction to Data Science / 24

Let's get started with Python basics! / 25

How to launch an interactive Python using IDLE	/ 26
Other online shells for Python	/ 27
Introduction to Data Types	/ 28
Let's learn how to code in Python	/ 29
Integers and Floats	/ 30
Basic operators	/ 32
Lesser-Known operators	/ 35
Variables	/ 35
Assignment operators	/ 37
Numbers & variables in the wild	/ 38
★ <i>Magic trick exercise</i>	/ 40
★ <i>Exercises for basic variables in Python</i>	/ 41

Strings Basic / 42

String operators	/ 43
String Indexing	/ 45
String Slices	/ 46
Print() function	/ 47
Escape characters	/ 47
Triple quotes	/ 49
More about strings	/ 51
★ <i>Exercises with strings</i>	/ 52

Introducing functions / 53

Len() function	/ 55
Input	/ 56
Type Casting	/ 58
F Strings	/ 60
★ <i>Age calculator exercise</i>	/ 60
★ <i>Shopping cart exercise</i>	/ 61

The world of methods	/ 63
Introducing methods Upper and Lower	/ 65
Method of navigating documentation in Python	/ 66
Help() & IPython	/ 68
Reading function Signatures + Strip methods	/ 68
Replace()	/ 70
Other very useful string methods for data researchers	/ 71
Method Chaining	/ 71
★ <i>Exercises with string methods</i>	/ 73
Booleans	/ 73
Comparison operators	/ 75
Comparing across types	/ 77
Truthiness & Falseyness	/ 78
The “in” operator	/ 79
Comparing strings	/ 80
★ <i>Exercises with Booleans</i>	/ 83
Conditionals basics	/ 84
Name length codealong	/ 86
A tangent on indentation	/ 87
Nesting conditionals	/ 88
★ <i>Water boiling Codealong</i>	/ 89
★ <i>BMI calculator exercise</i>	/ 89
★ <i>Tweet checker exercise</i>	/ 90
Writing more complex logic	/ 91
Logical AND	/ 92
Logical OR	/ 93
Logical NOT	/ 94
TruthyFalsey testing	/ 95
Logical operator precedence	/ 96
★ <i>Exercises with logical AND, OR, and NOT</i>	/ 97
Loops	/ 98
Avoiding infinite loops	/ 100
The range() function	/ 101
Working with Nested Loops	/ 102
Break and continue keywords	/ 102
★ <i>99 Bottles of Beer Codealong</i>	/ 104
★ <i>Loops problem set</i>	/ 104
★ <i>Snake Eyes Codealong</i>	/ 105
★ <i>Dice Roller Exercise</i>	/ 107
Functions	/ 108
Our very first function!	/ 109
Functions with an Input	/ 110
Functions with multiple arguments	/ 112
Introducing Return!	/ 113

- Using the Return keyword / 113
- Default parameters / 114
- Ordering default parameters / 115
- KeywordNamed argument / 116
- ★ *Function practice set* / 117
- Global Scope** / 118
 - Local Scope / 119
 - Scope in loops and conditionals / 120
 - Enclosing Scope / 120
 - Built-in Scope / 121
 - Scope precedence rules / 122
 - The 'Global' keyword / 123
 - ★ *Exercises for understanding the Global Scope* / 124
- Lists the basics** / 124
 - Accessing data in lists / 126
 - Updating list elements / 127
 - Append() and Extend() / 127
 - Insert() / 128
 - List Slices / 129
 - Deletion methods pop(), popitems(), remove() / 130
 - Iterating over lists / 131
 - Lists + loops patterns / 133
 - Nested lists / 135
 - List operators / 136
 - Sort(), Reverse(), and Count() / 137
 - Lists are mutable / 138
 - Comparing lists == vs is / 139
 - Join() and Split() / 139
 - List unpacking / 140
 - Copying lists / 141
 - ★ *Exercises with lists* / 143
 - ★ *Todo list exercise intro* / 144
- Dictionaries** / 145
 - Creating your Dictionaries / 146
 - Accessing data in Dictionaries / 148
 - Adding and updating data in Dictionaries / 149
 - The Get() method and "in" operator / 150
 - Dictionary Pop(), Clear(), and Del() / 151
 - Dictionaries are mutable too! / 152
 - Iterating Dicts Keys(), Values(), and Items() / 153
 - Fancy Dictionary merging / 154
 - Lists and Dicts combined / 155
 - Fromkeys() / 156
 - Update() / 157
 - ★ *Peak Dictionary exercise* / 158

Sets and Tuples / 159

- Tuple functionality / 161
- Sets introduction / 162
- Set operators: Intersection, Union, Difference / 163
- ★ *Exercises with sets* / 165

Back to functions. Introducing args / 165

- Introducing Kwargs / 167
- Parameter list ordering / 168
- A common gotcha mutable default Args / 169
- Unpacking Args / 170
- ★ *ArgsKwargs Exercises* / 170

Working with Errors / 171

- Common error types / 172
- Raising exceptions / 173
- When to raise / 174
- Try and except / 175
- LBYL and EAFP / 176
- ★ *Exercises with correct error handling* / 177

Modules / 178

- Working with built-in modules / 179
- Most popular built-in modules for Data Science / 180
- Fancy import syntax / 181
- Creating custom modules / 182
- 3rd party modules Pip & PyPI / 183
- Our first Pip package! / 183
- ★ *Sentiment analysis fun project installation* / 184

Object-Oriented Programming / 185

- Class Syntax / 186
- Writing our first class / 187
- Instance methods / 188
- ★ *Practicing Instance methods* / 189
- Class Attributes / 190
- Class Methods / 191
- Inheritance basics / 192
- The Super() function / 192

Part II: VISUALIZATION IN DATA SCIENCE**USING THE MOST IMPORTANT PYTHON MODULES** / 195**Introduction to Pandas module** / 197

- How to install Pandas? / 198
- Create a Series in Pandas / 199
- Create a DataFrame in Pandas / 200
- Read a CSV file with Pandas / 202

Advanced parameters /	203
Selecting rows and columns in Pandas /	204
Data wrangling in Pandas /	206
Arithmetics and statistics in Pandas /	210
Hierarchical indexing in Pandas /	212
Aggregation in Pandas /	214
Data Export in Pandas /	215
Pivot and Pivot Table in Pandas /	216
Visualization in Pandas /	217
★ <i>A few exercises with Pandas</i> /	229
NumPy, a perfect tool for working with Arrays /	230
Matrix Manipulation in Numpy /	237
Array Mathematics in Numpy /	239
Array Manipulation /	244
★ <i>Exercises with NumPy</i> /	251
Let's delve into DataVisualization in Python /	252
DataVisualization with Matplotlib /	253
★ <i>Matplotlib exercises for DataVisualization</i> /	256
Seaborn /	257
★ <i>Some exercises with Seaborn</i> /	265
Web Scraping in Data Science /	266
Scraping websites with Selenium /	272
★ <i>A few exercises with Webscraping</i> /	275
★ <i>Example of Gaussian noise (standard deviation)</i> /	275
Part III: INTRODUCTION TO BUSINESS STATISTICS IN DATA SCIENCE /	279
Big Data, Statistics, and Probability /	280
Business Intelligence (BI) techniques /	281
Big Data and Statistics /	282
Hypothesis testing /	285
Basic probability with Python /	286
Probability in Data Science /	288
Fundamentals of Combinatorics /	289
Bayes' Law /	291
Fundamentals of Probability Distributions /	292
★ <i>A Practical Example of Combinatorics</i> /	294
★ <i>A Practical Example of Bayesian Inference</i> /	295
Descriptive statistics /	297
Statistics with population and sample /	299
Cross Tables and Scatter Plots /	300
Skewness exercise solution /	302
★ <i>Exercises with Histograms in Descriptive Statistics</i> /	303
★ <i>Correlation exercise</i> /	303

Inferential Statistics	/	304
Inferential Statistics Confidence Intervals	/	306
The Normal Distribution	/	308
★ <i>Exercises with practical examples of Inferential Statistics</i>	/	310
Correlation and Regression	/	312
More about Correlation vs Regression	/	314
★ <i>Exercises with Correlation and Regression</i>	/	315
Time Series Analysis	/	316
Time Series Forecasting	/	317
Time Series – Visualization Basics	/	320
Time Series – Power Transformation	/	321
★ <i>Exercises with Time Series Analysis</i>	/	322
Part VI: MACHINE LEARNING (optional)	/	323
Scikit-learn, a free machine learning for advanced	/	323
Description of the start-up process	/	325
Description of the start-up process	/	325
Training and Test Data in Scikit-learn	/	327
Processing The Data Standardization in Scikit-learn	/	329
Normalizer class in Scikit-learn’s	/	329
Binarization in Scikit-learn	/	330
Encoding Categorical Features in Scikit-learn	/	331
Imputing Missing Values in Scikit-learn	/	333
Generating Polynomial Features in Scikit-learn	/	334
Create your model in Scikit-learn	/	335
Model Fitting in Scikit-learn	/	338
Prediction in Scikit-learn	/	339
Evaluate your model’s performance in Scikit-learn	/	340
Tune your model in Scikit-learn	/	341
★ <i>Exercises with Scikit-learn</i>	/	343
★ <i>How machine learning helps us as Data Scientists</i>	/	344
★ <i>Conclusions and tips for the future Data Science specialist</i>	/	346
★ <i>Where and what online materials</i>		
<i>we could read to learn more about Data Science</i>	/	347
★ <i>What kind of jobs can I find in</i>		
<i>Communication and Data Science, where, and how?</i>	/	348
★ <i>The impact of Data Science and</i>		
<i>Communication on the future of human society</i>	/	350
<i>References</i>	/	353

FOREWORD

Asking yourself the most important questions is a good way to introduce a book dedicated to Data Science because it helps to set the stage for the material that will be covered in the book. By asking questions, you can get a sense of what the book will be about and what you can expect to learn from it. Additionally, asking questions can help to engage the reader and get them thinking about the topic at hand, which can make the material more interesting and relevant to their own experiences. So, the first legitimate question is why is Data Science today a discipline of great importance for the present and future of master students? For sure, Data Science is a discipline of great importance today because it allows organizations to make better decisions by leveraging the vast amounts of data that are generated in today's world. With the help of Data Science, organizations can gain insights into their operations and customers, and use that information to improve their products, services, and overall business strategies. Additionally, Data Science is a rapidly growing field, with many job opportunities for individuals with the right skills and training. As a result, pursuing a master's degree in communication and Data Science can open a wide range of career possibilities for students.

Another key question is why should communication students be the ones to start and deepen Data Science or what is the connection between communication and public relations students and Data Science? Communication and public relations students should be interested in studying Data Science because it can help them better understand and analyze the vast amounts of data that are generated in today's world. Data Science can provide communication and public relations students with the skills to glean meaningful insights from data, which can be utilized to elevate and refine their practices in these areas. Learning Data Science, students can also gain the ability to extract valuable insights from data that can be applied to enhance their work in communication and public relations. Additionally, Data Sci-

ence can help communication and public relations students to better understand the needs and preferences of their target audience, and to create more effective communication strategies and campaigns. Finally, the skills and knowledge gained from studying Data Science can be highly valuable in the job market and can help communication and public relations students stand out from their peers and advance their careers.

A thirty-concrete key question would be how can it open a career and what are the shortcuts to success for any student wishing to specialize in this field? Studying Data Science can open a wide range of career possibilities for students. Data Science professionals may find employment as data scientists, data analysts, machine learning engineers, or business intelligence analysts. Alternatively, those with a strong background in Data Science may choose to pursue careers as data scientists, data analysts, machine learning engineers, or business intelligence analysts. In these roles, individuals can work in a variety of industries, including technology, finance, healthcare, and government, to help organizations make better decisions using data. To succeed in this field, students should be prepared to learn a wide range of technical skills, such as programming, statistics, and machine learning, as well as soft skills, such as problem-solving and communication. Additionally, students can gain a competitive advantage by participating in internships or other hands-on learning experiences, as well as by staying up to date with the latest developments in the field.

However, the field is very broad and then we should ask ourselves which parts or structures of Data Science are worth learning in the early stages. In this sense, what is worth studying in a master's program of only 2 years, having this limited time, and what is worth didactically deepened in the early stages of initiation? In a master's program with a limited time frame, students need to prioritize the key concepts and skills that are most essential for success in Data Science. Some of the most important areas to focus on in the early stages of a Data Science program include:

- In **Part I** of this book, we will learn **Fundamental Programming in Python: Data Science** involves working with large and complex

datasets, and students will need to be proficient in **Python**, to manipulate and analyze that data.

- In **Part II** of this book, we will learn about **Data Visualization** and other key procedures for a communication and data science specialist. Data visualization is an important tool for communicating the results of data analyses, and students will need to learn how to create clear and effective visualizations to effectively communicate their findings. **Web scraping** is an important tool for **data scientists** because it allows them to **extract data from websites** and turn it into structured, usable data that can be **analyzed** and **visualized**.
- In **Part III: Introduction to business statistics in Data Science**. This part aims to provide a fundamental understanding of statistics and its role in data science and business, as well as equip students with the tools and skills to apply statistical analysis to real-world situations.
- **Part IV** of this book is optional and is intended for readers who wish to further their knowledge in communication and data science. It covers the topic of **Machine learning** and is geared towards those who want to rapidly progress in this field and delve into advanced areas of data science. For the basics of this discipline, I recommend delving into just the first two parts, and optionally part three.

By focusing on these key areas, students can gain a solid foundation in Data Science, which they can then build upon in more advanced coursework and real-world experiences. The four stages of study broadly make up the overall structure of this book. By learning programming, statistics, machine learning, and data visualization, students can gain a solid foundation in Data Science, which they can then build upon in more advanced coursework and real-world experiences. These four areas are essential for success in Data Science because they provide the tools and techniques needed to manipulate, analyze, and interpret data, as well as to communicate the results of those analyses to others.

As the author of this book, I agree that organizing information into theoretical, example and practical components is an effective way to help students assimilate key information quickly and efficiently. By providing a mix of conceptual and practical material, students can

gain a deep understanding of the subject matter and apply this knowledge to real-world problems and scenarios. In addition, by providing examples and exercises, students can see how the concepts they are learning apply in practice, which can help to reinforce their understanding and facilitate the transfer of knowledge to new situations. Overall, this approach can help to engage and motivate learners and help them develop the skills and knowledge they need to succeed as communication and data scientists.

I hereby wish my students and all my readers to use this book to the fullest and learn to love Data Science professionally and then teach others as I have lovingly taught them. The best way to learn about any subject is to approach it with an open mind and a willingness to try new things. To maximize the benefit of this book, I would encourage my students to engage with the material actively, asking questions, trying out the examples and exercises, and seeking out additional resources and opportunities to learn more. Additionally, I encourage them to connect with other Data Science enthusiasts and professionals, either through online communities or in-person events, to learn from each other and stay up to date with the latest developments in the field. By adopting this approach, my students can develop a deep understanding of Data Science and they can become skilled and knowledgeable professionals who are well-equipped to teach others.