

Beginning Data Science in R

Data Analysis, Visualization,
and Modelling for the Data Scientist

Thomas Mailund

Apress

Beginning Data Science in R: Data Analysis, Visualization, and Modelling for the Data Scientist

Thomas Mailund
Aarhus, Denmark

ISBN-13 (pbk): 978-1-4842-2670-4
DOI 10.1007/978-1-4842-2671-1

ISBN-13 (electronic): 978-1-4842-2671-1

Library of Congress Control Number: 2017934529

Copyright © 2017 by Thomas Mailund

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/9781484226704. For more detailed information, please visit <http://www.apress.com/source-code>.

Contents at a Glance

Introduction	xxiii
■ Chapter 1: Introduction to R Programming	1
■ Chapter 2: Reproducible Analysis	29
■ Chapter 3: Data Manipulation	45
■ Chapter 4: Visualizing Data	75
■ Chapter 5: Working with Large Datasets	113
■ Chapter 6: Supervised Learning	125
■ Chapter 7: Unsupervised Learning	169
■ Chapter 8: More R Programming	205
■ Chapter 9: Advanced R Programming	233
■ Chapter 10: Object Oriented Programming	257
■ Chapter 11: Building an R Package	269
■ Chapter 12: Testing and Package Checking	281
■ Chapter 13: Version Control	287
■ Chapter 14: Profiling and Optimizing	303
Index	347

Contents

Introduction	xxiii
■ Chapter 1: Introduction to R Programming	1
Basic Interaction with R	1
Using R as a Calculator	3
Simple Expressions	3
Assignments	5
Actually, All of the Above Are Vectors of Values.	5
Indexing Vectors	6
Vectorized Expressions.....	7
Comments	8
Functions.....	8
Getting Documentation for Functions	9
Writing Your Own Functions.....	10
Vectorized Expressions and Functions	12
A Quick Look at Control Structures	12
Factors	16
Data Frames	18
Dealing with Missing Values.....	20
Using R Packages.....	21

Data Pipelines (or Pointless Programming)	22
Writing Pipelines of Function Calls	23
Writing Functions that Work with Pipelines	23
The magical “.” argument	24
Defining Functions Using .	25
Anonymous Functions	26
Other Pipeline Operations	27
Coding and Naming Conventions	28
Exercises	28
Mean of Positive Values	28
Root Mean Square Error	28
Chapter 2: Reproducible Analysis	29
Literate Programming and Integration of Workflow and Documentation	30
Creating an R Markdown/knitr Document in RStudio	30
The YAML Language	33
The Markdown Language	34
Formatting Text	35
Cross-Referencing	38
Bibliographies	39
Controlling the Output (Templates/Stylesheets)	39
Running R Code in Markdown Documents	40
Using Chunks when Analyzing Data (Without Compiling Documents)	42
Caching Results	43
Displaying Data	43
Exercises	44
Create an R Markdown Document	44
Produce Different Output	44
Add Caching	44

Chapter 3: Data Manipulation	45
Data Already in R.....	45
Quickly Reviewing Data.....	47
Reading Data.....	48
Examples of Reading and Formatting Datasets	49
Breast Cancer Dataset.....	49
Boston Housing Dataset	55
The readr Package.....	56
Manipulating Data with dplyr	58
Some Useful dplyr Functions.....	59
Breast Cancer Data Manipulation.....	65
Tidying Data with tidyr	69
Exercises	72
Importing Data.....	73
Using dplyr.....	73
Using tidyr	73
Chapter 4: Visualizing Data	75
Basic Graphics.....	75
The Grammar of Graphics and the ggplot2 Package.....	83
Using qplot()	84
Using Geometries	88
Facets	97
Scaling.....	100
Themes and Other Graphics Transformations.....	105
Figures with Multiple Plots.....	109
Exercises	111

Chapter 5: Working with Large Datasets	113
Subsample Your Data Before You Analyze the Full Dataset.....	113
Running Out of Memory During Analysis.....	115
Too Large to Plot.....	116
Too Slow to Analyze.....	120
Too Large to Load.....	121
Exercises.....	124
Subsampling.....	124
Hex and 2D Density Plots.....	124
Chapter 6: Supervised Learning	125
Machine Learning.....	125
Supervised Learning.....	125
Regression versus Classification.....	126
Inference versus Prediction.....	127
Specifying Models.....	128
Linear Regression.....	128
Logistic Regression (Classification, Really).....	133
Model Matrices and Formula.....	136
Validating Models.....	145
Evaluating Regression Models.....	145
Evaluating Classification Models.....	147
Random Permutations of Your Data.....	153
Cross-Validation.....	157
Selecting Random Training and Testing Data.....	159
Examples of Supervised Learning Packages.....	161
Decision Trees.....	161
Random Forests.....	163
Neural Networks.....	164
Support Vector Machines.....	165

Naive Bayes.....	165
Exercises	166
Fitting Polynomials	166
Evaluating Different Classification Measures	166
Breast Cancer Classification.....	166
Leave-One-Out Cross-Validation (Slightly More Difficult).....	167
Decision Trees	167
Random Forests.....	167
Neural Networks.....	167
Support Vector Machines.....	167
Compare Classification Algorithms.....	167
■ Chapter 7: Unsupervised Learning	169
Dimensionality Reduction.....	169
Principal Component Analysis	169
Multidimensional Scaling	177
Clustering	181
k-Means Clustering	182
Hierarchical Clustering	188
Association Rules	192
Exercises	196
Dealing with Missing Data in the HouseVotes84 Data.....	196
Rescaling for k-Means Clustering	196
Varying k.....	196
Project 1	196
Importing Data.....	197
Exploring the Data	198
Fitting Models.....	203

Exercises	204
Exploring Other Formulas	204
Exploring Different Models	204
Analyzing Your Own Dataset.....	204
■ Chapter 8: More R Programming.....	205
Expressions	205
Arithmetic Expressions	205
Boolean Expressions	206
Basic Data Types	207
The Numeric Type	207
The Integer Type	208
The Complex Type.....	208
The Logical Type	208
The Character Type	209
Data Structures	209
Vectors.....	209
Matrix	210
Lists	212
Indexing.....	213
Named Values.....	215
Factors.....	216
Formulas.....	216
Control Structures	216
Selection Statements	216
Loops	218
A Word of Warning About Looping	219
Functions.....	220
Named Arguments	221
Default Parameters.....	222
Return Values.....	222

Lazy Evaluation.....	223
Scoping.....	224
Function Names Are Different from Variable Names	227
Recursive Functions	227
Exercises	229
Fibonacci Numbers.....	229
Outer Product	229
Linear Time Merge.....	229
Binary Search	230
More Sorting.....	230
Selecting the k Smallest Element.....	231
■ Chapter 9: Advanced R Programming	233
Working with Vectors and Vectorizing Functions	233
ifelse	235
Vectorizing Functions	235
The apply Family	237
Advanced Functions	242
Special Names.....	242
Infix Operators	242
Replacement Functions.....	243
How Mutable Is Data Anyway?	245
Functional Programming.....	246
Anonymous Functions	246
Functions Taking Functions as Arguments	247
Functions Returning Functions (and Closures).....	247
Filter, Map, and Reduce	248
Function Operations: Functions as Input and Output	250
Ellipsis Parameters.....	253

Exercises	255
between.....	255
apply_if.....	255
power.....	255
Row and Column Sums	255
Factorial Again.....	255
Function Composition	256
■ Chapter 10: Object Oriented Programming	257
Immutable Objects and Polymorphic Functions.....	257
Data Structures	257
Example: Bayesian Linear Model Fitting.....	258
Classes	259
Polymorphic Functions.....	261
Defining Your Own Polymorphic Functions.....	262
Class Hierarchies.....	263
Specialization as Interface	263
Specialization in Implementations.....	264
Exercises	267
Shapes.....	267
Polynomials	267
■ Chapter 11: Building an R Package	269
Creating an R Package	269
Package Names.....	269
The Structure of an R Package.....	270
.Rbuildignore	270
Description	271
NAMESPACE.....	274
R/ and man/	275

Roxygen.....	275
Documenting Functions.....	275
Import and Export.....	276
Package Scope Versus Global Scope.....	277
Internal Functions.....	277
File Load Order	277
Adding Data to Your Package	278
Building an R Package	279
Exercises	280
■ Chapter 12: Testing and Package Checking	281
Unit Testing.....	281
Automating Testing.....	282
Using testthat	283
Writing Good Tests	284
Using Random Numbers in Tests.....	285
Testing Random Results	285
Checking a Package for Consistency	286
Exercise.....	286
■ Chapter 13: Version Control.....	287
Version Control and Repositories	287
Using git in RStudio.....	288
Installing git.....	288
Making Changes to Files, Staging Files, and Committing Changes.....	289
Adding git to an Existing Project.....	291
Bare Repositories and Cloning Repositories.....	291
Pushing Local Changes and Fetching and Pulling Remote Changes.....	292
Handling Conflicts.....	294
Working with Branches	294
Typical Workflows Involve Lots of Branches.....	297
Pushing Branches to the Global Repository.....	297

GitHub.....	297
Moving an Existing Repository to GitHub.....	299
Installing Packages from GitHub	300
Collaborating on GitHub.....	300
Pull Requests.....	300
Forking Repositories Instead of Cloning.....	301
Exercises	301
■ Chapter 14: Profiling and Optimizing	303
Profiling	303
A Graph-Flow Algorithm	304
Speeding Up Your Code	315
Parallel Execution.....	317
Switching to C++	320
Exercises	322
Project 2	322
Bayesian Linear Regression	323
Exercises: Priors and Posteriors.....	324
Predicting Target Variables for New Predictor Values.....	328
Formulas and Their Model Matrix.....	330
Working with Model Matrices in R.....	331
Exercises	334
Model Matrices Without Response Variables.....	334
Exercises	335
Interface to a blm Class	336
Constructor.....	336
Updating Distributions: An Example Interface	337
Designing Your blm Class	340
Model Methods.....	340

Building an R Package for blm	342
Deciding on the Package Interface.....	342
Organization of Source Files.....	342
Document Your Package Interface Well.....	343
Adding README and NEWS Files to Your Package	343
Testing.....	344
GitHub.....	344
Conclusions	344
Data Science.....	345
Machine Learning.....	345
Data Analysis	345
R Programming.....	345
The End	346
Acknowledgements.....	346
Index.....	347

Introduction

Welcome to *Introduction to Data Science with R*. This book was written as a set of lecture notes for two classes I teach, *Data Science: Visualization and Analysis* and *Data Science: Software Development and Testing*. The book is written to fit the structure of these classes, where each class consists of seven weeks of lectures and project work. This means that there are 14 chapters with the core material, where the first seven focus on data analysis and the last seven on developing reusable software for data science.

What Is Data Science?

Oh boy! That is a difficult question. I don't know if it is easy to find someone who is entirely sure what data science is, but I am pretty sure that it would be difficult to find two people with fewer than three opinions about it. It is certainly a popular buzzword, and everyone wants to have data scientists these days, so data science skills are useful to have on the CV. But what *is* it?

Since I can't really give you an agreed upon definition, I will just give you my own: *Data science is the science of learning from data.*

This is a very broad definition—almost too broad to be useful. I realize this. But then, I think data science is an incredibly general field. I don't have a problem with that. Of course, you could argue that any *science* is all about getting information out of data, and you might be right. Although I would say that there is more to science than just transforming raw data into useful information. The sciences are focusing on answering specific questions about the world while data science is focusing on how to manipulate data efficiently and effectively. The primary focus is not which questions to ask of the data but how we can answer them, whatever they may be. It is more like computer science and mathematics than it is like natural sciences, in this way. It isn't so much about studying the natural world as it is about how to compute data efficiently.

Included in data science is the design of experiments. With the right data, we can address the questions we are interested in. With a poor design of experiments or a poor choice of which data we gather, this can be difficult. Study design might be the most important aspect of data science, but is not the topic of this book. In this book I focus on the analysis of data, once gathered.

Computer science is also mainly the study of computations—as is hinted at in the name—but is a bit broader in this focus. Although *datalogy*, an earlier name for data science, was also suggested for computer science, and for example in Denmark it *is* the name for computer science, using the name “computer science” puts the focus on computation while using the name “data science” puts the focus on data. But of course, the fields overlap. If you are writing a sorting algorithm, are you then focusing on the computation or the data? Is that even a meaningful question to ask?

There is a huge overlap between computer science and data science and naturally the skillsets you need overlap as well. To efficiently manipulate data you need the tools for doing that, so computer programming skills are a must and some knowledge about algorithms and data structures usually is as well. For data science, though, the focus is always on the data. In a data analysis project, the focus is on how the data flows from its raw form through various manipulations until it is summarized in some useful form. Although the difference can be subtle, the focus is not about what operations a program does during the analysis, but about how the data flows and is transformed. It is also focused on *why* we do certain transformations of the

data, what purpose those changes serve, and how they help us gain knowledge about the data. It is as much about deciding what to do with the data as it is about how to do it efficiently.

Statistics is of course also closely related to data science. So closely linked, in fact, that many consider data science just a fancy word for statistics that looks slightly more modern and sexy. I can't say that I strongly disagree with this—data science *does* sound sexier than statistics—but just as data science is slightly different from computer science, data science is also slightly different from statistics. Just, perhaps, somewhat less different than computer science is.

A large part of doing statistics is building mathematical models for your data and fitting the models to the data to learn about the data in this way. That is also what we do in data science. As long as the focus is on the data, I am happy to call statistics data science. If the focus changes to the models and the mathematics, then we are drifting away from data science into something else—just as if the focus changes from the data to computations we are drifting from data science to computer science.

Data science is also related to machine learning and artificial intelligence, and again there are huge overlaps. Perhaps not surprising since something like machine learning has its home both in computer science and in statistics; if it is focusing on data analysis, it is also at home in data science. To be honest, it has never been clear to me when a mathematical model changes from being a plain old statistical model to becoming machine learning anyway.

For this book, we are just going to go with my definition and, as long as we are focusing on analyzing data, we are going to call it data science.

Prerequisites for Reading this Book

In the first seven chapters in this book, the focus is on data analysis and not programming. For those seven chapters, I do not assume a detailed familiarity with topics such as software design, algorithms, data structures, and such. I do not expect you to have any experience with the R programming language either. I do, however, expect that you have had *some* experience with programming, mathematical modeling, and statistics.

Programming R can be quite tricky at times if you are familiar with a scripting language or object-oriented languages. R is a functional language that does not allow you to modify data, and while it does have systems for object-oriented programming, it handles this programming paradigm very differently from languages you are likely to have seen such as Java or Python.

For the data analysis part of this book, the first seven chapters, we will only use R for very straightforward programming tasks, so none of this should pose a problem. We will have to write simple scripts for manipulating and summarizing data so you should be familiar with how to write basic expressions like function calls, `if` statements, loops, and so on. These things you will have to be comfortable with. I will introduce every such construction in the book when we need them so you will see how they are expressed in R, but I will not spend much time explaining them. I mostly will just expect you to be able to pick it up from examples.

Similarly, I do not expect you to know already how to fit data and compare models in R. I do expect that you have had enough introduction to statistics to be comfortable with basic terms like parameter estimation, model fitting, explanatory and response variables, and model comparison. If not, I expect you to be at least able to pick up what we are talking about when you need to.

I won't expect you to know a lot about statistics and programming, but this isn't *Data Science for Dummies*, so I do expect you to be able to figure out examples without me explaining everything in detail.

After the first seven chapters is a short description of a data analysis project, one of my students did in an earlier class. It shows how such a project could look, but I suggest that you do not wait until you have finished the first seven chapters to start doing such analysis yourself. To get the most benefit out of reading this book, you should be applying what you learn continuously. Already when you begin reading, I suggest that you find a dataset that you would be interested in finding out more about and then apply what you learn in each chapter to that data.

For the final seven chapters of the book, the focus *is* on programming. To read this part you should be familiar with object-oriented programming. I will explain how it is handled in R and how it differs from languages such as Python, Java or C++ but I expect you to be familiar with terms such as class hierarchies, inheritance, and polymorphic methods. I will not expect you to be already familiar with functional programming (but if you are, there should still be plenty to learn in those chapters if you are not already familiar with R programming as well).

Plan for the Book

In the book, we cover basic data manipulation—filtering and selecting relevant data; transforming data into shapes readily analyzable; summarizing data; visualizing data in informative ways both for exploring data and presenting results; and model building. These are the key aspects of doing analysis in data science. After this we will cover how to develop R code that is reusable and works well with existing packages, and that is easy to extend, and we will see how to build new R packages that other people will be able to use in their projects. These are the essential skills you will need to develop your own methods and share them with the world.

We will do all this using the programming language R (<https://www.r-project.org/about.html>). R is one of the most popular (and open source) data analysis programming languages around at the moment. Of course, popularity doesn't imply quality, but because R is so popular it has a rich ecosystem of extensions (called “packages” in R) for just about any kind of analysis you could be interested in. People who develop statistical methods often implement them as R packages, so you can quite often get the state of the art techniques very easily in R. The popularity also means that there is a large community of people who can help if you have problems. Most problems you run into can be solved with a few minutes on Google because you are unlikely to be the first to run into any particular issue. There are also plenty of online tutorials for learning more about R and specialized packages, there are plenty of videos with talks about R and popular R packages, and there are plenty of books you can buy if you want to learn more.

Data Analysis and Visualization

The topics focusing on data analysis and visualization are covered in the first seven chapters:

- Chapter 1, Introduction to R programming. In which you learn how to work with data and write data pipelines.
- Chapter 2, **Reproducible** analysis. In which you find out how to integrate documentation and analysis in a single document and how to use such documents to produce reproducible research.
- Chapter 3, Data manipulation. In which you learn how to import, tidy up, and transform data, and compute summaries from data.
- Chapter 4, **Visualizing** and exploring data. In which you learn how to make plots for exploring data features and for presenting data features and analysis results.
- Chapter 5, **Working** with large datasets. In which you learn how to deal with data where the number of observations make the usual approaches too slow.
- Chapter 6, **Supervised** learning. In which you learn how to train models when you have datasets with known classes or regression values.
- Chapter 7, **Unsupervised** learning. In which you learn how to search for patterns you are not aware of in data.

These chapters are followed by the first project, where you see the various techniques in use.

Software Development

Software and package development is then covered in the following seven chapters:

- Chapter 8, **More** R programming. In which you'll return to the basics of R programming and get a few more details than the tutorial in Chapter 1.
- Chapter 9, **Advanced R programming**. In which you explore more advanced features of the R programming language, in particular, functional programming.
- Chapter 10, **Object** oriented programming. In which you learn how R models object orientation and how you can use it to write more generic code.
- Chapter 11, **Building** an R package. In which you learn the necessary components of an R package and how to program your own.
- Chapter 12, **Testing** and checking. In which you learn techniques for testing your R code and checking the consistency of your R packages.
- Chapter 13, **Version** control. In which you learn how to manage code under version control and how to collaborate using GitHub.
- Chapter 14, **Profiling** and optimizing. In which you learn how to identify hotspots of code where inefficient solutions are slowing you down and techniques for alleviating this.

These chapters are then followed by the second project, where you'll build a package for Bayesian linear regression.

Getting R and RStudio

You will need to install R on your computer to do the exercises in this book. I suggest that you get an integrated environment since it can be slightly easier to keep track of a project when you have your plots, documentation, code, etc., all in the same program.

I personally use RStudio (<https://www.rstudio.com/products/RStudio>), which I warmly recommend. You can get it for free—just follow the link—and I will assume that you have it when I need to refer to the software environment you are using in the following chapters. There won't be much RStudio specifics, though, and most tools for working with R have the same features, so if you want to use something else you can probably follow the notes without any difficulties.

Projects

You cannot learn how to analyze data without analyzing data, and you cannot learn how to develop software without developing software either. Typing in examples from the book is nothing like writing code on your own. Even doing exercises from the book—which you really ought to do—is not the same as working on your own projects. Exercises, after all, cover small isolated aspects of problems you have just been introduced to. In the real world, there is not a chapter of material presented before every task you have to deal with. You need to work out by yourself what needs to be done and how. If you only do the exercises in this book, you will miss the most important lessons in analyzing data. How to explore the data and get a feeling for it; how to do the detective work necessary to pull out some understanding from the data; and how to deal with all the noise and weirdness found in any dataset. And for developing a package, you need to think through how to design and implement its functionality so that the various functions and data structures fit well together.

In this book, I go through a data analysis project to show you what that can look like. To actually learn how to analyze data, you need to do it yourself as well, and you need to do it with a dataset that I haven't analyzed for you. You might have a dataset lying around you have worked on before, a dataset from something you are just interested in, or you can probably find something interesting at a public data repository, e.g., one of these:

- RDataMining.com
- [UCI machine learning repository \(http://archive.ics.uci.edu/ml/\)](http://archive.ics.uci.edu/ml/)
- [KDNuggets \(http://www.kdnuggets.com/datasets/index.html\)](http://www.kdnuggets.com/datasets/index.html)
- [Reddit r/datasets \(https://www.reddit.com/r/datasets\)](https://www.reddit.com/r/datasets)
- [GitHub awesome public datasets \(https://github.com/caesar0301/awesome-public-datasets\)](https://github.com/caesar0301/awesome-public-datasets)

I suggest that you find yourself a dataset and that after each lesson, you use the skills you have learned to explore this dataset. Pick data that is structured as a table with observations as rows and variables as columns, since that is the form of the data we consider in this book. At the end of the first seven chapters, you will have analyzed this data, you can write a report about your analysis that others can evaluate to follow and maybe modify it. You will be doing reproducible science.

For the programming topics, I describe another project illustrating the design and implementation issues involved in making an R package. There, you should be able to learn from just implementing your own version of the project I use, but you will, of course, be more challenged by working on a project without any of my help at all. Whichever you do, to get the full benefit of this book you should make your own package while reading the programming chapters.

Introduction to R Programming

We will use R for our data analysis so we need to know the basics of programming in the R language. R is a full programming language with both functional programming and object oriented programming features. Learning the language is far beyond the scope of this chapter and is something we return to later. The good news, though, is that to use R for data analysis, you rarely need to do much programming. At least, if you do the *right* kind of programming, you won't need much.

For manipulating data—and how to do this is the topic of the next chapter—you mainly just have to string together a couple of operations. Operations such as “group the data by this feature” followed by “calculate the mean value of these features within each group” and then “plot these means”. This used to be much more complicated to do in R, but a couple of new ideas on how to structure such data flow—and some clever implementations of these in a couple of packages such as `magrittr` and `dplyr`—has significantly simplified it. We will see some of this at the end of this chapter and more in the next chapter. First, though, you need to get a taste for R.

Basic Interaction with R

Start by downloading RStudio if you haven't done so already (<https://www.rstudio.com/products/RStudio>). If you open it, you should see a window similar to Figure 1-1. Well, except that you will be in an empty project while the figure shows (on the top right) that this RStudio is opened in a project called “Data Science”. You always want to be working on a project. Projects keep track of the state of your analysis by remembering variables and functions you have written and keep track of which files you have opened and such. Choose File ► New Project to create a project. You can create a project from an existing directory, but if this is the first time you are working with R you probably just want to create an empty project in a new directory, so do that.