

Data Visualization with R

Rob Kabacoff

2018-09-03

Contents

| | |
|--|------------|
| Welcome | 7 |
| Preface | 9 |
| How to use this book | 9 |
| Prerequisites | 10 |
| Setup | 10 |
| 1 Data Preparation | 11 |
| 1.1 Importing data | 11 |
| 1.2 Cleaning data | 12 |
| 2 Introduction to ggplot2 | 19 |
| 2.1 A worked example | 19 |
| 2.2 Placing the data and mapping options | 30 |
| 2.3 Graphs as objects | 32 |
| 3 Univariate Graphs | 35 |
| 3.1 Categorical | 35 |
| 3.2 Quantitative | 51 |
| 4 Bivariate Graphs | 63 |
| 4.1 Categorical vs. Categorical | 63 |
| 4.2 Quantitative vs. Quantitative | 71 |
| 4.3 Categorical vs. Quantitative | 79 |
| 5 Multivariate Graphs | 103 |
| 5.1 Grouping | 103 |
| 6 Maps | 115 |
| 6.1 Dot density maps | 115 |
| 6.2 Choropleth maps | 119 |

| | | |
|-----------|-------------------------------|------------|
| 7 | Time-dependent graphs | 127 |
| 7.1 | Time series | 127 |
| 7.2 | Dummbbell charts | 130 |
| 7.3 | Slope graphs | 133 |
| 7.4 | Area Charts | 135 |
| 8 | Statistical Models | 139 |
| 8.1 | Correlation plots | 139 |
| 8.2 | Linear Regression | 141 |
| 8.3 | Logistic regression | 145 |
| 8.4 | Survival plots | 147 |
| 8.5 | Mosaic plots | 150 |
| 9 | Other Graphs | 153 |
| 9.1 | 3-D Scatterplot | 153 |
| 9.2 | Biplots | 159 |
| 9.3 | Bubble charts | 161 |
| 9.4 | Flow diagrams | 163 |
| 9.5 | Heatmaps | 168 |
| 9.6 | Radar charts | 174 |
| 9.7 | Scatterplot matrix | 176 |
| 9.8 | Waterfall charts | 178 |
| 9.9 | Word clouds | 180 |
| 10 | Customizing Graphs | 183 |
| 10.1 | Axes | 183 |
| 10.2 | Colors | 187 |
| 10.3 | Points & Lines | 193 |
| 10.4 | Legends | 195 |
| 10.5 | Labels | 197 |
| 10.6 | Annotations | 199 |
| 10.7 | Themes | 206 |
| 11 | Saving Graphs | 219 |
| 11.1 | Via menus | 219 |
| 11.2 | Via code | 219 |
| 11.3 | File formats | 219 |
| 11.4 | External editing | 221 |

| | |
|--|------------|
| 12 Interactive Graphs | 223 |
| 12.1 leaflet | 223 |
| 12.2 plotly | 223 |
| 12.3 rbokeh | 226 |
| 12.4 rCharts | 226 |
| 12.5 highcharter | 226 |
| 13 Advice / Best Practices | 231 |
| 13.1 Labeling | 231 |
| 13.2 Signal to noise ratio | 232 |
| 13.3 Color choice | 234 |
| 13.4 y -Axis scaling | 234 |
| 13.5 Attribution | 238 |
| 13.6 Going further | 238 |
| 13.7 Final Note | 239 |
| A Datasets | 241 |
| A.1 Academic salaries | 241 |
| A.2 Starwars | 241 |
| A.3 Mammal sleep | 241 |
| A.4 Marriage records | 242 |
| A.5 Fuel economy data | 242 |
| A.6 Gapminder data | 242 |
| A.7 Current Population Survey (1985) | 242 |
| A.8 Houston crime data | 242 |
| A.9 US economic timeseries | 243 |
| A.10 Saratoga housing data | 243 |
| A.11 US population by age and year | 243 |
| A.12 NCCTG lung cancer data | 243 |
| A.13 Titanic data | 243 |
| A.14 JFK Cuban Missile speech | 244 |
| A.15 UK Energy forecast data | 244 |
| A.16 US Mexican American Population | 244 |
| B About the Author | 245 |
| C About the QAC | 247 |

Welcome

R is an amazing platform for data analysis, capable of creating almost any type of graph. This book helps you create the most popular visualizations - from quick and dirty plots to publication-ready graphs. The text relies heavily on the ggplot2 package for graphics, but other approaches are covered as well.

This work is licensed under a Creative Commons Attribution-NonCommercial-NoDerivatives 4.0 International License.

My goal is make this book as helpful and user-friendly as possible. Any feedback is both welcome and appreciated.

Preface

How to use this book

You don't need to read this book from start to finish in order to start building effective graphs. Feel free to jump to the section that you need and then explore others that you find interesting.

Graphs are organized by

- the number of variables to be plotted
- the type of variables to be plotted
- the purpose of the visualization

| Chapter | Description |
|-----------------------|---|
| Ch 1 | provides a quick overview of how to get your data into R and how to prepare it for analysis. |
| Ch 2 | provides an overview of the <code>ggplot2</code> package. |
| Ch 3 | describes graphs for visualizing the distribution of a single categorical (e.g. race) or quantitative (e.g. income) variable. |
| Ch 4 | describes graphs that display the relationship between two variables. |
| Ch 5 | describes graphs that display the relationships among 3 or more variables. It is helpful to read chapters 3 and 4 before this chapter. |
| Ch 6 | provides a brief introduction to displaying data geographically. |
| Ch 7 | describes graphs that display change over time. |
| Ch 8 | describes graphs that can help you interpret the results of statistical models. |
| Ch 9 | covers graphs that do not fit neatly elsewhere (every book needs a miscellaneous chapter). |
| Ch 10 | describes how to customize the look and feel of your graphs. If you are going to share your graphs with others, be sure to skim this chapter. |
| Ch 11 | covers how to save your graphs. Different formats are optimized for different purposes. |
| Ch 12 | provides an introduction to interactive graphics. |
| Ch 13 | gives advice on creating effective graphs and where to go to learn more. It's worth a look. |
| The Appendices | describe each of the datasets used in this book, and provides a short blurb about the author and the Wesleyan Quantitative Analysis Center. |

There is **no one right graph** for displaying data. Check out the examples, and see which type best fits your needs.

Prerequisites

It's assumed that you have some experience with the R language and that you have already installed R and RStudio. If not, here are some resources for getting started:

- A (very) short introduction to R
- DataCamp - Introduction to R with Jonathon Cornelissen
- Quick-R
- Getting up to speed with R

Setup

In order to create the graphs in this guide, you'll need to install some optional R packages. To install **all** of the necessary packages, run the following code in the RStudio console window.

```
pkgs <- c("ggplot2", "dplyr", "tidyr",
          "mosaicData", "carData",
          "VIM", "scales", "treemapify",
          "gapminder", "ggmap", "choroplethr",
          "choroplethrMaps", "CGPfunctions",
          "ggcorrplot", "visreg",
          "gcookbook", "forcats",
          "survival", "survminer",
          "ggalluvial", "ggridges",
          "GGally", "superheat",
          "waterfalls", "factoextra",
          "networkD3", "ggthemes",
          "hrbrthemes", "ggpol",
          "ggbeeswarm")
install.packages(pkgs)
```

Alternatively, you can install a given package the first time it is needed.

For example, if you execute

```
library(gapminder)
```

and get the message

```
Error in library(gapminder) : there is no package called 'gapminder'
```

you know that the package has never been installed. Simply execute

```
install.packages("gapminder")
```

once and

```
library(gapminder)
```

will work from that point on.