

DATA SCIENCE SERIES

# INTRODUCTION TO DATA SCIENCE

Data Analysis and Prediction  
Algorithms with R

RAFAEL A. IRIZARRY



CRC Press  
Taylor & Francis Group

A CHAPMAN & HALL BOOK

# Introduction to Data Science

Data Analysis and Prediction  
Algorithms with R

## CHAPMAN & HALL/CRC DATA SCIENCE SERIES

Reflecting the interdisciplinary nature of the field, this book series brings together researchers, practitioners, and instructors from statistics, computer science, machine learning, and analytics. The series will publish cutting-edge research, industry applications, and textbooks in data science.

The inclusion of concrete examples, applications, and methods is highly encouraged. The scope of the series includes titles in the areas of machine learning, pattern recognition, predictive analytics, business analytics, Big Data, visualization, programming, software, learning analytics, data wrangling, interactive graphics, and reproducible research.

Published Titles

### **Probability and Statistics for Data Science: Math + R + Data**

Norman Matloff

### **Feature Engineering and Selection: A Practical Approach for Predictive Models**

Max Kuhn and Kjell Johnson

### **Introduction to Data Science: Data Analysis and Prediction Algorithms with R**

Rafael A. Irizarry

# Introduction to Data Science

Data Analysis and Prediction  
Algorithms with R

Rafael A. Irizarry

Dana-Farber Cancer Institute and Harvard University



CRC Press

Taylor & Francis Group  
Boca Raton London New York

---

CRC Press is an imprint of the  
Taylor & Francis Group, an **informa** business

CRC Press  
Taylor & Francis Group  
6000 Broken Sound Parkway NW, Suite 300  
Boca Raton, FL 33487-2742

© 2020 by Taylor & Francis Group, LLC  
CRC Press is an imprint of Taylor & Francis Group, an Informa business

No claim to original U.S. Government works

Printed on acid-free paper

International Standard Book Number-13: 978-0-367-35798-6 (Hardback)

This book contains information obtained from authentic and highly regarded sources. Reasonable efforts have been made to publish reliable data and information, but the author and publisher cannot assume responsibility for the validity of all materials or the consequences of their use. The authors and publishers have attempted to trace the copyright holders of all material reproduced in this publication and apologize to copyright holders if permission to publish in this form has not been obtained. If any copyright material has not been acknowledged please write and let us know so we may rectify in any future reprint.

Except as permitted under U.S. Copyright Law, no part of this book may be reprinted, reproduced, transmitted, or utilized in any form by any electronic, mechanical, or other means, now known or hereafter invented, including photocopying, microfilming, and recording, or in any information storage or retrieval system, without written permission from the publishers.

For permission to photocopy or use material electronically from this work, please access [www.copyright.com](http://www.copyright.com) (<http://www.copyright.com/>) or contact the Copyright Clearance Center, Inc. (CCC), 222 Rosewood Drive, Danvers, MA 01923, 978-750-8400. CCC is a not-for-profit organization that provides licenses and registration for a variety of users. For organizations that have been granted a photocopy license by the CCC, a separate system of payment has been arranged.

**Trademark Notice:** Product or corporate names may be trademarks or registered trademarks, and are used only for identification and explanation without intent to infringe.

---

**Library of Congress Cataloging-in-Publication Data**

---

Names: Irizarry, Rafael A., author.  
Title: Introduction to data science : data analysis and prediction  
algorithms with R / Rafael A. Irizarry.  
Description: [Boca Raton] : [CRC Press], [2019] | Summary: "The book begins by going over the basics of R and the tidyverse. You learn R throughout the book, but in the first part we go over the building blocks needed to keep learning during the rest of the book"-- Provided by publisher.  
Identifiers: LCCN 2019025160 (print) | LCCN 2019025161 (ebook) | ISBN 9780367357986 (hardback) | ISBN 9780367357993 (paperback) | ISBN 9780429341830 (Adobe PDF)  
Subjects: LCSH: R (Computer program language) | Information visualization. | Data mining. | Statistics--Data processing. | Probabilities--Data processing. | Computer algorithms. | Quantitative research.  
Classification: LCC QA276.45.R3 I75 2019 (print) | LCC QA276.45.R3 (ebook) | DDC 005.362--dc23  
LC record available at <https://lcn.loc.gov/2019025160>  
LC ebook record available at <https://lcn.loc.gov/2019025161>

---

Visit the Taylor & Francis Web site at  
<http://www.taylorandfrancis.com>

and the CRC Press Web site at  
<http://www.crcpress.com>

---

# Contents

---

<b>Preface</b>	<b>xxv</b>
<b>Acknowledgments</b>	<b>xxvii</b>
<b>Introduction</b>	<b>xxix</b>
<b>1 Getting started with R and RStudio</b>	<b>1</b>
1.1 Why R? . . . . .	1
1.2 The R console . . . . .	1
1.3 Scripts . . . . .	2
1.4 RStudio . . . . .	3
1.4.1 The panes . . . . .	3
1.4.2 Key bindings . . . . .	5
1.4.3 Running commands while editing scripts . . . . .	6
1.4.4 Changing global options . . . . .	8
1.5 Installing R packages . . . . .	8
<b>I R</b>	<b>11</b>
<b>2 R basics</b>	<b>13</b>
2.1 Case study: US Gun Murders . . . . .	13
2.2 The very basics . . . . .	15
2.2.1 Objects . . . . .	15
2.2.2 The workspace . . . . .	15
2.2.3 Functions . . . . .	16
2.2.4 Other prebuilt objects . . . . .	18
2.2.5 Variable names . . . . .	18
2.2.6 Saving your workspace . . . . .	19
2.2.7 Motivating scripts . . . . .	19
2.2.8 Commenting your code . . . . .	20
2.3 Exercises . . . . .	20

2.4	Data types . . . . .	21
2.4.1	Data frames . . . . .	21
2.4.2	Examining an object . . . . .	21
2.4.3	The accessor: <code>\$</code> . . . . .	22
2.4.4	Vectors: numerics, characters, and logical . . . . .	23
2.4.5	Factors . . . . .	24
2.4.6	Lists . . . . .	24
2.4.7	Matrices . . . . .	26
2.5	Exercises . . . . .	27
2.6	Vectors . . . . .	28
2.6.1	Creating vectors . . . . .	28
2.6.2	Names . . . . .	29
2.6.3	Sequences . . . . .	29
2.6.4	Subsetting . . . . .	30
2.7	Coercion . . . . .	31
2.7.1	Not availables (NA) . . . . .	32
2.8	Exercises . . . . .	32
2.9	Sorting . . . . .	33
2.9.1	<code>sort</code> . . . . .	33
2.9.2	<code>order</code> . . . . .	33
2.9.3	<code>max</code> and <code>which.max</code> . . . . .	34
2.9.4	<code>rank</code> . . . . .	34
2.9.5	Beware of recycling . . . . .	35
2.10	Exercises . . . . .	35
2.11	Vector arithmetics . . . . .	36
2.11.1	Rescaling a vector . . . . .	37
2.11.2	Two vectors . . . . .	37
2.12	Exercises . . . . .	38
2.13	Indexing . . . . .	38
2.13.1	Subsetting with logicals . . . . .	38
2.13.2	Logical operators . . . . .	39
2.13.3	<code>which</code> . . . . .	39
2.13.4	<code>match</code> . . . . .	40
2.13.5	<code>%in%</code> . . . . .	40

2.14	Exercises	40
2.15	Basic plots	41
2.15.1	<code>plot</code>	41
2.15.2	<code>hist</code>	42
2.15.3	<code>boxplot</code>	43
2.15.4	<code>image</code>	43
2.16	Exercises	44
<b>3</b>	<b>Programming basics</b>	<b>45</b>
3.1	Conditional expressions	45
3.2	Defining functions	47
3.3	Namespaces	48
3.4	For-loops	49
3.5	Vectorization and functionals	50
3.6	Exercises	51
<b>4</b>	<b>The tidyverse</b>	<b>53</b>
4.1	Tidy data	53
4.2	Exercises	54
4.3	Manipulating data frames	55
4.3.1	Adding a column with <code>mutate</code>	55
4.3.2	Subsetting with <code>filter</code>	56
4.3.3	Selecting columns with <code>select</code>	56
4.4	Exercises	57
4.5	The pipe: <code>%&gt;%</code>	58
4.6	Exercises	59
4.7	Summarizing data	60
4.7.1	<code>summarize</code>	60
4.7.2	<code>pull</code>	62
4.7.3	Group then summarize with <code>group_by</code>	63
4.8	Sorting data frames	64
4.8.1	Nested sorting	64
4.8.2	The top $n$	65
4.9	Exercises	65
4.10	Tibbles	66



4.10.1	Tibbles display better . . . . .	67
4.10.2	Subsets of tibbles are tibbles . . . . .	67
4.10.3	Tibbles can have complex entries . . . . .	68
4.10.4	Tibbles can be grouped . . . . .	68
4.10.5	Create a tibble using <code>tibble</code> instead of <code>data.frame</code> . . . . .	68
4.11	The dot operator . . . . .	69
4.12	<code>do</code> . . . . .	70
4.13	The <b>purrr</b> package . . . . .	71
4.14	Tidyverse conditionals . . . . .	73
4.14.1	<code>case_when</code> . . . . .	73
4.14.2	<code>between</code> . . . . .	73
4.15	Exercises . . . . .	74
<b>5</b>	<b>Importing data</b>	<b>75</b>
5.1	Paths and the working directory . . . . .	76
5.1.1	The filesystem . . . . .	76
5.1.2	Relative and full paths . . . . .	77
5.1.3	The working directory . . . . .	77
5.1.4	Generating path names . . . . .	78
5.1.5	Copying files using paths . . . . .	78
5.2	The readr and readxl packages . . . . .	79
5.2.1	readr . . . . .	79
5.2.2	readxl . . . . .	80
5.3	Exercises . . . . .	80
5.4	Downloading files . . . . .	81
5.5	R-base importing functions . . . . .	82
5.5.1	<code>scan</code> . . . . .	82
5.6	Text versus binary files . . . . .	83
5.7	Unicode versus ASCII . . . . .	83
5.8	Organizing data with spreadsheets . . . . .	84
5.9	Exercises . . . . .	84

<b>II</b>	<b>Data Visualization</b>	<b>85</b>
<b>6</b>	<b>Introduction to data visualization</b>	<b>87</b>
<b>7</b>	<b>ggplot2</b>	<b>91</b>
7.1	The components of a graph . . . . .	92
7.2	ggplot objects . . . . .	93
7.3	Geometries . . . . .	94
7.4	Aesthetic mappings . . . . .	95
7.5	Layers . . . . .	96
7.5.1	Tinkering with arguments . . . . .	97
7.6	Global versus local aesthetic mappings . . . . .	98
7.7	Scales . . . . .	99
7.8	Labels and titles . . . . .	100
7.9	Categories as colors . . . . .	101
7.10	Annotation, shapes, and adjustments . . . . .	102
7.11	Add-on packages . . . . .	103
7.12	Putting it all together . . . . .	104
7.13	Quick plots with <code>qplot</code> . . . . .	105
7.14	Grids of plots . . . . .	106
7.15	Exercises . . . . .	106
<b>8</b>	<b>Visualizing data distributions</b>	<b>109</b>
8.1	Variable types . . . . .	109
8.2	Case study: describing student heights . . . . .	110
8.3	Distribution function . . . . .	110
8.4	Cumulative distribution functions . . . . .	111
8.5	Histograms . . . . .	112
8.6	Smoothed density . . . . .	113
8.6.1	Interpreting the y-axis . . . . .	117
8.6.2	Densities permit stratification . . . . .	118
8.7	Exercises . . . . .	118
8.8	The normal distribution . . . . .	122
8.9	Standard units . . . . .	124
8.10	Quantile-quantile plots . . . . .	125
8.11	Percentiles . . . . .	127

8.12	Boxplots . . . . .	127
8.13	Stratification . . . . .	129
8.14	Case study: describing student heights (continued) . . . . .	129
8.15	Exercises . . . . .	131
8.16	ggplot2 geometries . . . . .	132
8.16.1	Barplots . . . . .	133
8.16.2	Histograms . . . . .	134
8.16.3	Density plots . . . . .	135
8.16.4	Boxplots . . . . .	136
8.16.5	QQ-plots . . . . .	136
8.16.6	Images . . . . .	137
8.16.7	Quick plots . . . . .	138
8.17	Exercises . . . . .	140
<b>9</b>	<b>Data visualization in practice</b>	<b>141</b>
9.1	Case study: new insights on poverty . . . . .	141
9.1.1	Hans Rosling's quiz . . . . .	142
9.2	Scatterplots . . . . .	143
9.3	Faceting . . . . .	144
9.3.1	<code>facet_wrap</code> . . . . .	146
9.3.2	Fixed scales for better comparisons . . . . .	147
9.4	Time series plots . . . . .	147
9.4.1	Labels instead of legends . . . . .	150
9.5	Data transformations . . . . .	151
9.5.1	Log transformation . . . . .	151
9.5.2	Which base? . . . . .	153
9.5.3	Transform the values or the scale? . . . . .	154
9.6	Visualizing multimodal distributions . . . . .	155
9.7	Comparing multiple distributions with boxplots and ridge plots . . . . .	155
9.7.1	Boxplots . . . . .	156
9.7.2	Ridge plots . . . . .	157
9.7.3	Example: 1970 versus 2010 income distributions . . . . .	159
9.7.4	Accessing computed variables . . . . .	164
9.7.5	Weighted densities . . . . .	167
9.8	The ecological fallacy and importance of showing the data . . . . .	167

9.8.1	Logistic transformation . . . . .	168
9.8.2	Show the data . . . . .	168
<b>10</b>	<b>Data visualization principles</b>	<b>171</b>
10.1	Encoding data using visual cues . . . . .	171
10.2	Know when to include 0 . . . . .	174
10.3	Do not distort quantities . . . . .	177
10.4	Order categories by a meaningful value . . . . .	179
10.5	Show the data . . . . .	180
10.6	Ease comparisons . . . . .	183
10.6.1	Use common axes . . . . .	183
10.6.2	Align plots vertically to see horizontal changes and horizontally to see vertical changes . . . . .	184
10.6.3	Consider transformations . . . . .	185
10.6.4	Visual cues to be compared should be adjacent . . . . .	187
10.6.5	Use color . . . . .	188
10.7	Think of the color blind . . . . .	188
10.8	Plots for two variables . . . . .	189
10.8.1	Slope charts . . . . .	189
10.8.2	Bland-Altman plot . . . . .	191
10.9	Encoding a third variable . . . . .	191
10.10	Avoid pseudo-three-dimensional plots . . . . .	193
10.11	Avoid too many significant digits . . . . .	195
10.12	Know your audience . . . . .	196
10.13	Exercises . . . . .	196
10.14	Case study: vaccines and infectious diseases . . . . .	201
10.15	Exercises . . . . .	204
<b>11</b>	<b>Robust summaries</b>	<b>205</b>
11.1	Outliers . . . . .	205
11.2	Median . . . . .	206
11.3	The inter quartile range (IQR) . . . . .	206
11.4	Tukey's definition of an outlier . . . . .	207
11.5	Median absolute deviation . . . . .	208
11.6	Exercises . . . . .	208
11.7	Case study: self-reported student heights . . . . .	209

<b>III Statistics with R</b>	<b>213</b>
<b>12 Introduction to statistics with R</b>	<b>215</b>
<b>13 Probability</b>	<b>217</b>
13.1 Discrete probability . . . . .	217
13.1.1 Relative frequency . . . . .	217
13.1.2 Notation . . . . .	218
13.1.3 Probability distributions . . . . .	218
13.2 Monte Carlo simulations for categorical data . . . . .	218
13.2.1 Setting the random seed . . . . .	220
13.2.2 With and without replacement . . . . .	220
13.3 Independence . . . . .	221
13.4 Conditional probabilities . . . . .	221
13.5 Addition and multiplication rules . . . . .	222
13.5.1 Multiplication rule . . . . .	222
13.5.2 Multiplication rule under independence . . . . .	222
13.5.3 Addition rule . . . . .	223
13.6 Combinations and permutations . . . . .	223
13.6.1 Monte Carlo example . . . . .	227
13.7 Examples . . . . .	227
13.7.1 Monty Hall problem . . . . .	228
13.7.2 Birthday problem . . . . .	229
13.8 Infinity in practice . . . . .	231
13.9 Exercises . . . . .	232
13.10 Continuous probability . . . . .	234
13.11 Theoretical continuous distributions . . . . .	235
13.11.1 Theoretical distributions as approximations . . . . .	235
13.11.2 The probability density . . . . .	237
13.12 Monte Carlo simulations for continuous variables . . . . .	238
13.13 Continuous distributions . . . . .	239
13.14 Exercises . . . . .	239
<b>14 Random variables</b>	<b>241</b>
14.1 Random variables . . . . .	241
14.2 Sampling models . . . . .	242

14.3	The probability distribution of a random variable . . . . .	243
14.4	Distributions versus probability distributions . . . . .	245
14.5	Notation for random variables . . . . .	245
14.6	The expected value and standard error . . . . .	246
14.6.1	Population SD versus the sample SD . . . . .	248
14.7	Central Limit Theorem . . . . .	249
14.7.1	How large is large in the Central Limit Theorem? . . . . .	250
14.8	Statistical properties of averages . . . . .	250
14.9	Law of large numbers . . . . .	252
14.9.1	Misinterpreting law of averages . . . . .	252
14.10	Exercises . . . . .	252
14.11	Case study: The Big Short . . . . .	254
14.11.1	Interest rates explained with chance model . . . . .	254
14.11.2	The Big Short . . . . .	257
14.12	Exercises . . . . .	260
<b>15</b>	<b>Statistical inference</b>	<b>261</b>
15.1	Polls . . . . .	261
15.1.1	The sampling model for polls . . . . .	262
15.2	Populations, samples, parameters, and estimates . . . . .	264
15.2.1	The sample average . . . . .	264
15.2.2	Parameters . . . . .	265
15.2.3	Polling versus forecasting . . . . .	265
15.2.4	Properties of our estimate: expected value and standard error . . . . .	266
15.3	Exercises . . . . .	267
15.4	Central Limit Theorem in practice . . . . .	268
15.4.1	A Monte Carlo simulation . . . . .	269
15.4.2	The spread . . . . .	271
15.4.3	Bias: why not run a very large poll? . . . . .	271
15.5	Exercises . . . . .	272
15.6	Confidence intervals . . . . .	274
15.6.1	A Monte Carlo simulation . . . . .	276
15.6.2	The correct language . . . . .	277
15.7	Exercises . . . . .	277
15.8	Power . . . . .	278

15.9	p-values . . . . .	279
15.10	Association tests . . . . .	280
15.10.1	Lady Tasting Tea . . . . .	281
15.10.2	Two-by-two tables . . . . .	282
15.10.3	Chi-square Test . . . . .	282
15.10.4	The odds ratio . . . . .	283
15.10.5	Confidence intervals for the odds ratio . . . . .	284
15.10.6	Small count correction . . . . .	285
15.10.7	Large samples, small p-values . . . . .	285
15.11	Exercises . . . . .	286
<b>16</b>	<b>Statistical models</b>	<b>287</b>
16.1	Poll aggregators . . . . .	288
16.1.1	Poll data . . . . .	290
16.1.2	Pollster bias . . . . .	292
16.2	Data-driven models . . . . .	293
16.3	Exercises . . . . .	295
16.4	Bayesian statistics . . . . .	298
16.4.1	Bayes theorem . . . . .	298
16.5	Bayes theorem simulation . . . . .	299
16.5.1	Bayes in practice . . . . .	300
16.6	Hierarchical models . . . . .	301
16.7	Exercises . . . . .	303
16.8	Case study: election forecasting . . . . .	305
16.8.1	Bayesian approach . . . . .	306
16.8.2	The general bias . . . . .	307
16.8.3	Mathematical representations of models . . . . .	307
16.8.4	Predicting the electoral college . . . . .	310
16.8.5	Forecasting . . . . .	314
16.9	Exercises . . . . .	317
16.10	The t-distribution . . . . .	318
<b>17</b>	<b>Regression</b>	<b>321</b>
17.1	Case study: is height hereditary? . . . . .	321
17.2	The correlation coefficient . . . . .	322

17.2.1	Sample correlation is a random variable . . . . .	324
17.2.2	Correlation is not always a useful summary . . . . .	326
17.3	Conditional expectations . . . . .	326
17.4	The regression line . . . . .	329
17.4.1	Regression improves precision . . . . .	330
17.4.2	Bivariate normal distribution (advanced) . . . . .	331
17.4.3	Variance explained . . . . .	333
17.4.4	Warning: there are two regression lines . . . . .	333
17.5	Exercises . . . . .	334
<b>18</b>	<b>Linear models</b>	<b>335</b>
18.1	Case study: Moneyball . . . . .	335
18.1.1	Sabermetrics . . . . .	336
18.1.2	Baseball basics . . . . .	337
18.1.3	No awards for BB . . . . .	338
18.1.4	Base on balls or stolen bases? . . . . .	339
18.1.5	Regression applied to baseball statistics . . . . .	341
18.2	Confounding . . . . .	344
18.2.1	Understanding confounding through stratification . . . . .	345
18.2.2	Multivariate regression . . . . .	348
18.3	Least squares estimates . . . . .	348
18.3.1	Interpreting linear models . . . . .	349
18.3.2	Least Squares Estimates (LSE) . . . . .	349
18.3.3	The <code>lm</code> function . . . . .	351
18.3.4	LSE are random variables . . . . .	352
18.3.5	Predicted values are random variables . . . . .	353
18.4	Exercises . . . . .	354
18.5	Linear regression in the tidyverse . . . . .	355
18.5.1	The broom package . . . . .	358
18.6	Exercises . . . . .	359
18.7	Case study: Moneyball (continued) . . . . .	360
18.7.1	Adding salary and position information . . . . .	364
18.7.2	Picking nine players . . . . .	365
18.8	The regression fallacy . . . . .	367
18.9	Measurement error models . . . . .	369



18.10 Exercises . . . . .	371
<b>19 Association is not causation</b>	<b>373</b>
19.1 Spurious correlation . . . . .	373
19.2 Outliers . . . . .	376
19.3 Reversing cause and effect . . . . .	378
19.4 Confounders . . . . .	379
19.4.1 Example: UC Berkeley admissions . . . . .	379
19.4.2 Confounding explained graphically . . . . .	380
19.4.3 Average after stratifying . . . . .	381
19.5 Simpson's paradox . . . . .	382
19.6 Exercises . . . . .	383
<b>IV Data Wrangling</b>	<b>385</b>
<b>20 Introduction to data wrangling</b>	<b>387</b>
<b>21 Reshaping data</b>	<b>389</b>
21.1 <code>gather</code> . . . . .	389
21.2 <code>spread</code> . . . . .	391
21.3 <code>separate</code> . . . . .	391
21.4 <code>unite</code> . . . . .	394
21.5 Exercises . . . . .	395
<b>22 Joining tables</b>	<b>397</b>
22.1 Joins . . . . .	398
22.1.1 Left join . . . . .	399
22.1.2 Right join . . . . .	400
22.1.3 Inner join . . . . .	400
22.1.4 Full join . . . . .	400
22.1.5 Semi join . . . . .	401
22.1.6 Anti join . . . . .	401
22.2 Binding . . . . .	402
22.2.1 Binding columns . . . . .	402
22.2.2 Binding by rows . . . . .	402
22.3 Set operators . . . . .	403
22.3.1 Intersect . . . . .	403

22.3.2	Union . . . . .	404
22.3.3	<code>setdiff</code> . . . . .	404
22.3.4	<code>setequal</code> . . . . .	404
22.4	Exercises . . . . .	405
<b>23</b>	<b>Web scraping</b>	<b>407</b>
23.1	HTML . . . . .	408
23.2	The <code>rvest</code> package . . . . .	409
23.3	CSS selectors . . . . .	411
23.4	JSON . . . . .	412
23.5	Exercises . . . . .	413
<b>24</b>	<b>String processing</b>	<b>415</b>
24.1	The <code>stringr</code> package . . . . .	415
24.2	Case study 1: US murders data . . . . .	417
24.3	Case study 2: self-reported heights . . . . .	419
24.4	How to <i>escape</i> when defining strings . . . . .	421
24.5	Regular expressions . . . . .	423
24.5.1	Strings are a regexp . . . . .	423
24.5.2	Special characters . . . . .	423
24.5.3	Character classes . . . . .	425
24.5.4	Anchors . . . . .	426
24.5.5	Quantifiers . . . . .	426
24.5.6	White space <code>\s</code> . . . . .	427
24.5.7	Quantifiers: <code>*</code> , <code>?</code> , <code>+</code> . . . . .	428
24.5.8	Not . . . . .	428
24.5.9	Groups . . . . .	429
24.6	Search and replace with regex . . . . .	430
24.6.1	Search and replace using groups . . . . .	432
24.7	Testing and improving . . . . .	433
24.8	Trimming . . . . .	435
24.9	Changing lettercase . . . . .	436
24.10	Case study 2: self-reported heights (continued) . . . . .	436
24.10.1	The <code>extract</code> function . . . . .	437
24.10.2	Putting it all together . . . . .	438

24.11	String splitting . . . . .	439
24.12	Case study 3: extracting tables from a PDF . . . . .	442
24.13	Recoding . . . . .	445
24.14	Exercises . . . . .	446
<b>25</b>	<b>Parsing dates and times</b>	<b>449</b>
25.1	The date data type . . . . .	449
25.2	The lubridate package . . . . .	450
25.3	Exercises . . . . .	453
<b>26</b>	<b>Text mining</b>	<b>455</b>
26.1	Case study: Trump tweets . . . . .	455
26.2	Text as data . . . . .	457
26.3	Sentiment analysis . . . . .	462
26.4	Exercises . . . . .	467
<b>V</b>	<b>Machine Learning</b>	<b>469</b>
<b>27</b>	<b>Introduction to machine learning</b>	<b>471</b>
27.1	Notation . . . . .	471
27.2	An example . . . . .	472
27.3	Exercises . . . . .	474
27.4	Evaluation metrics . . . . .	474
27.4.1	Training and test sets . . . . .	475
27.4.2	Overall accuracy . . . . .	476
27.4.3	The confusion matrix . . . . .	478
27.4.4	Sensitivity and specificity . . . . .	479
27.4.5	Balanced accuracy and $F_1$ score . . . . .	481
27.4.6	Prevalence matters in practice . . . . .	482
27.4.7	ROC and precision-recall curves . . . . .	483
27.4.8	The loss function . . . . .	484
27.5	Exercises . . . . .	486
27.6	Conditional probabilities and expectations . . . . .	486
27.6.1	Conditional probabilities . . . . .	487
27.6.2	Conditional expectations . . . . .	488
27.6.3	Conditional expectation minimizes squared loss function . . . . .	488
27.7	Exercises . . . . .	489

27.8	Case study: is it a 2 or a 7? . . . . .	489
<b>28</b>	<b>Smoothing</b>	<b>493</b>
28.1	Bin smoothing . . . . .	495
28.2	Kernels . . . . .	497
28.3	Local weighted regression (loess) . . . . .	498
28.3.1	Fitting parabolas . . . . .	502
28.3.2	Beware of default smoothing parameters . . . . .	503
28.4	Connecting smoothing to machine learning . . . . .	504
28.5	Exercises . . . . .	504
<b>29</b>	<b>Cross validation</b>	<b>507</b>
29.1	Motivation with k-nearest neighbors . . . . .	507
29.1.1	Over-training . . . . .	509
29.1.2	Over-smoothing . . . . .	510
29.1.3	Picking the $k$ in kNN . . . . .	511
29.2	Mathematical description of cross validation . . . . .	513
29.3	K-fold cross validation . . . . .	514
29.4	Exercises . . . . .	517
29.5	Bootstrap . . . . .	518
29.6	Exercises . . . . .	521
<b>30</b>	<b>The caret package</b>	<b>523</b>
30.1	The caret <code>train</code> function . . . . .	523
30.2	Cross validation . . . . .	524
30.3	Example: fitting with loess . . . . .	526
<b>31</b>	<b>Examples of algorithms</b>	<b>529</b>
31.1	Linear regression . . . . .	529
31.1.1	The <code>predict</code> function . . . . .	530
31.2	Exercises . . . . .	531
31.3	Logistic regression . . . . .	533
31.3.1	Generalized linear models . . . . .	534
31.3.2	Logistic regression with more than one predictor . . . . .	538
31.4	Exercises . . . . .	539
31.5	k-nearest neighbors . . . . .	540
31.6	Exercises . . . . .	541

31.7	Generative models . . . . .	541
31.7.1	Naive Bayes . . . . .	542
31.7.2	Controlling prevalence . . . . .	543
31.7.3	Quadratic discriminant analysis . . . . .	545
31.7.4	Linear discriminant analysis . . . . .	547
31.7.5	Connection to distance . . . . .	549
31.8	Case study: more than three classes . . . . .	549
31.9	Exercises . . . . .	553
31.10	Classification and regression trees (CART) . . . . .	554
31.10.1	The curse of dimensionality . . . . .	554
31.10.2	CART motivation . . . . .	555
31.10.3	Regression trees . . . . .	558
31.10.4	Classification (decision) trees . . . . .	564
31.11	Random forests . . . . .	566
31.12	Exercises . . . . .	571
<b>32</b>	<b>Machine learning in practice</b>	<b>573</b>
32.1	Preprocessing . . . . .	574
32.2	k-nearest neighbor and random forest . . . . .	575
32.3	Variable importance . . . . .	578
32.4	Visual assessments . . . . .	579
32.5	Ensembles . . . . .	579
32.6	Exercises . . . . .	580
<b>33</b>	<b>Large datasets</b>	<b>581</b>
33.1	Matrix algebra . . . . .	581
33.1.1	Notation . . . . .	582
33.1.2	Converting a vector to a matrix . . . . .	584
33.1.3	Row and column summaries . . . . .	585
33.1.4	<code>apply</code> . . . . .	586
33.1.5	Filtering columns based on summaries . . . . .	586
33.1.6	Indexing with matrices . . . . .	588
33.1.7	Binarizing the data . . . . .	590
33.1.8	Vectorization for matrices . . . . .	590
33.1.9	Matrix algebra operations . . . . .	591

<i>Contents</i>	xxi
33.2 Exercises . . . . .	591
33.3 Distance . . . . .	591
33.3.1 Euclidean distance . . . . .	592
33.3.2 Distance in higher dimensions . . . . .	592
33.3.3 Euclidean distance example . . . . .	593
33.3.4 Predictor space . . . . .	595
33.3.5 Distance between predictors . . . . .	595
33.4 Exercises . . . . .	595
33.5 Dimension reduction . . . . .	596
33.5.1 Preserving distance . . . . .	596
33.5.2 Linear transformations (advanced) . . . . .	599
33.5.3 Orthogonal transformations (advanced) . . . . .	600
33.5.4 Principal component analysis . . . . .	602
33.5.5 Iris example . . . . .	604
33.5.6 MNIST example . . . . .	607
33.6 Exercises . . . . .	609
33.7 Recommendation systems . . . . .	610
33.7.1 Movielens data . . . . .	610
33.7.2 Recommendation systems as a machine learning challenge . . . . .	612
33.7.3 Loss function . . . . .	612
33.7.4 A first model . . . . .	613
33.7.5 Modeling movie effects . . . . .	614
33.7.6 User effects . . . . .	615
33.8 Exercises . . . . .	616
33.9 Regularization . . . . .	617
33.9.1 Motivation . . . . .	617
33.9.2 Penalized least squares . . . . .	619
33.9.3 Choosing the penalty terms . . . . .	622
33.10 Exercises . . . . .	624
33.11 Matrix factorization . . . . .	625
33.11.1 Factors analysis . . . . .	628
33.11.2 Connection to SVD and PCA . . . . .	630
33.12 Exercises . . . . .	633

<b>34 Clustering</b>	<b>639</b>
34.1 Hierarchical clustering . . . . .	640
34.2 k-means . . . . .	642
34.3 Heatmaps . . . . .	642
34.4 Filtering features . . . . .	643
34.5 Exercises . . . . .	644
<b>VI Productivity Tools</b>	<b>645</b>
<b>35 Introduction to productivity tools</b>	<b>647</b>
<b>36 Organizing with Unix</b>	<b>649</b>
36.1 Naming convention . . . . .	649
36.2 The terminal . . . . .	650
36.3 The filesystem . . . . .	650
36.3.1 Directories and subdirectories . . . . .	651
36.3.2 The home directory . . . . .	651
36.3.3 Working directory . . . . .	652
36.3.4 Paths . . . . .	653
36.4 Unix commands . . . . .	653
36.4.1 <code>ls</code> : Listing directory content . . . . .	654
36.4.2 <code>mkdir</code> and <code>rmdir</code> : make and remove a directory . . . . .	654
36.4.3 <code>cd</code> : navigating the filesystem by changing directories . . . . .	655
36.5 Some examples . . . . .	657
36.6 More Unix commands . . . . .	658
36.6.1 <code>mv</code> : moving files . . . . .	658
36.6.2 <code>cp</code> : copying files . . . . .	659
36.6.3 <code>rm</code> : removing files . . . . .	659
36.6.4 <code>less</code> : looking at a file . . . . .	659
36.7 Preparing for a data science project . . . . .	660
36.8 Advanced Unix . . . . .	661
36.8.1 Arguments . . . . .	661
36.8.2 Getting help . . . . .	662
36.8.3 Pipes . . . . .	662
36.8.4 Wild cards . . . . .	663
36.8.5 Environment variables . . . . .	663

36.8.6	Shells . . . . .	664
36.8.7	Executables . . . . .	664
36.8.8	Permissions and file types . . . . .	665
36.8.9	Commands you should learn . . . . .	665
36.8.10	File manipulation in R . . . . .	665
<b>37</b>	<b>Git and GitHub</b>	<b>667</b>
37.1	Why use Git and GitHub? . . . . .	667
37.2	GitHub accounts . . . . .	667
37.3	GitHub repositories . . . . .	670
37.4	Overview of Git . . . . .	671
37.4.1	Clone . . . . .	672
37.5	Initializing a Git directory . . . . .	676
37.6	Using Git and GitHub in RStudio . . . . .	678
<b>38</b>	<b>Reproducible projects with RStudio and R markdown</b>	<b>683</b>
38.1	RStudio projects . . . . .	683
38.2	R markdown . . . . .	686
38.2.1	The header . . . . .	688
38.2.2	R code chunks . . . . .	688
38.2.3	Global options . . . . .	689
38.2.4	knitr . . . . .	689
38.2.5	More on R markdown . . . . .	690
38.3	Organizing a data science project . . . . .	690
38.3.1	Create directories in Unix . . . . .	690
38.3.2	Create an RStudio project . . . . .	691
38.3.3	Edit some R scripts . . . . .	692
38.3.4	Create some more directories using Unix . . . . .	693
38.3.5	Add a README file . . . . .	693
38.3.6	Initializing a Git directory . . . . .	693
38.3.7	Add, commit, and push files using RStudio . . . . .	694
	<b>Index</b>	<b>695</b>





# Taylor & Francis

Taylor & Francis Group

<http://taylorandfrancis.com>

---

# *Preface*

---

This book started out as the class notes used in the HarvardX Data Science Series<sup>1</sup>.

The link for the online version of the book is <https://rafalab.github.io/dsbook/>

The R markdown code used to generate the book is available on GitHub<sup>2</sup>. Note that, the graphical theme used for plots throughout the book can be recreated using the `ds_theme_set()` function from **dslabs** package.

This work is licensed under the Creative Commons Attribution-NonCommercial-ShareAlike 4.0 International (CC BY-NC-SA 4.0)<sup>3</sup>.

We make announcements related to the book on Twitter. For updates follow @rafalab<sup>4</sup>

---

<sup>1</sup><https://www.edx.org/professional-certificate/harvardx-data-science>

<sup>2</sup><https://github.com/rafalab/dsbook>

<sup>3</sup><https://creativecommons.org/licenses/by-nc-sa/4.0>

<sup>4</sup><https://twitter.com/rafalab>



# Taylor & Francis

Taylor & Francis Group

<http://taylorandfrancis.com>

---

# Acknowledgments

---

This book is dedicated to all the people involved in building and maintaining R and the R packages we use in this book. A special thanks to the developers and maintainers of R base, the tidyverse, and the caret package.

A special thanks to my tidyverse guru David Robinson and Amy Gill for dozens of comments, edits, and suggestions. Also, many thanks to Stephanie Hicks who twice served as a co-instructor in my data science classes and Yihui Xie who patiently put up with my many questions about bookdown. Thanks also to Karl Broman, from whom I borrowed ideas for the Data Visualization and Productivity Tools parts, and to Hector Corrada-Bravo, for advice on how to best teach machine learning. Thanks to Peter Aldhous from whom I borrowed ideas for the principles of data visualization section and Jenny Bryan for writing *Happy Git and GitHub for the useR*, which influenced our Git chapters. Thanks to Alyssa Frazee for helping create the homework problem that became the Recommendation Systems chapter and to Amanda Cox for providing the New York Regents exams data. Also, many thanks to Jeff Leek, Roger Peng, and Brian Caffo, whose class inspired the way this book is divided and to Garrett Grolmund and Hadley Wickham for making the bookdown code for their R for Data Science book open. Finally, thanks to Alex Nones for proofreading the manuscript during its various stages.

This book was conceived during the teaching of several applied statistics courses, starting over fifteen years ago. The teaching assistants working with me throughout the years made important indirect contributions to this book. The latest iteration of this course is a HarvardX series coordinated by Heather Sternshein and Zzofia Gajdos. We thank them for their contributions. We are also grateful to all the students whose questions and comments helped us improve the book. The courses were partially funded by NIH grant R25GM114818. We are very grateful to the National Institutes of Health for its support.

A special thanks goes to all those who edited the book via GitHub pull requests or made suggestions by creating an *issue*: `nickyfoto` (Huang Qiang), `desautm` (Marc-André Désautels), `michaschwab` (Michail Schwab), `alvarolarreategui` (Alvaro Larreategui), `jakevc` (Jake VanCampen), `omerta` (Guillermo Lengemann), `espinielli` (Enrico Spinielli), `asimumba` (Aaron Simumba), `braunschweig` (Maldewar), `gwierzchowski` (Grzegorz Wierzchowski), `technocrat` (Richard Careaga), `atzakas`, `defeit` (David Emerson Feit), `shiraamitchell` (Shira Mitchell), `Nathalie-S`, `andreashandel` (Andreas Handel), `berkowitze` (Elias Berkowitz), `Dean-Webb` (Dean Webber), `mohayusuf`, `jinrothstein`, `mPloenzke` (Matthew Ploenzke), and David D. Kane.



# Taylor & Francis

Taylor & Francis Group

<http://taylorandfrancis.com>

---

# Introduction

---

The demand for skilled data science practitioners in industry, academia, and government is rapidly growing. This book introduces concepts and skills that can help you tackle real-world data analysis challenges. It covers concepts from probability, statistical inference, linear regression, and machine learning. It also helps you develop skills such as R programming, data wrangling with **dplyr**, data visualization with **ggplot2**, algorithm building with **caret**, file organization with UNIX/Linux shell, version control with Git and GitHub, and reproducible document preparation with **knitr** and R markdown. The book is divided into six parts: **R**, **Data Visualization**, **Statistics with R**, **Data Wrangling**, **Machine Learning**, and **Productivity Tools**. Each part has several chapters meant to be presented as one lecture and includes dozens of exercises distributed across chapters.

---

## Case studies

Throughout the book, we use motivating case studies. In each case study, we try to realistically mimic a data scientist’s experience. For each of the concepts covered, we start by asking specific questions and answer these through data analysis. We learn the concepts as a means to answer the questions. Examples of the case studies included in the book are:

Case Study	Concept
US murder rates by state	R Basics
Student heights	Statistical Summaries
Trends in world health and economics	Data Visualization
The impact of vaccines on infectious disease rates	Data Visualization
The financial crisis of 2007-2008	Probability
Election forecasting	Statistical Inference
Reported student heights	Data Wrangling
Money Ball: Building a baseball team	Linear Regression
MNIST: Image processing hand-written digits	Machine Learning
Movie recommendation systems	Machine Learning

---

## Who will find this book useful?

This book is meant to be a textbook for a first course in Data Science. No previous knowledge of R is necessary, although some experience with programming may be helpful. The statistical concepts used to answer the case study questions are only briefly introduced, so a Probability

and Statistics textbook is highly recommended for in-depth understanding of these concepts. If you read and understand all the chapters and complete all the exercises, you will be well-positioned to perform basic data analysis tasks and you will be prepared to learn the more advanced concepts and skills needed to become an expert.

---

## What does this book cover?

We start by going over the **basics of R** and the **tidyverse**. You learn R throughout the book, but in the first part we go over the building blocks needed to keep learning.

The growing availability of informative datasets and software tools has led to increased reliance on **data visualizations** in many fields. In the second part we demonstrate how to use **ggplot2** to generate graphs and describe important data visualization principles.

In the third part we demonstrate the importance of statistics in data analysis by answering case study questions using **probability, inference, and regression** with R.

The fourth part uses several examples to familiarize the reader with **data wrangling**. Among the specific skills we learn are web scrapping, using regular expressions, and joining and reshaping data tables. We do this using **tidyverse** tools.

In the fifth part we present several challenges that lead us to introduce **machine learning**. We learn to use the **caret** package to build prediction algorithms including K-nearest neighbors and random forests.

In the final part, we provide a brief introduction to the **productivity tools** we use on a day-to-day basis in data science projects. These are RStudio, UNIX/Linux shell, Git and GitHub, and **knitr** and R Markdown.

---

## What is not covered by this book?

This book focuses on the data analysis aspects of data science. We therefore do not cover aspects related to data management or engineering. Although R programming is an essential part of the book, we do not teach more advanced computer science topics such as data structures, optimization, and algorithm theory. Similarly, we do not cover topics such as web services, interactive graphics, parallel computing, and data streaming processing. The statistical concepts are presented mainly as tools to solve problems and in-depth theoretical descriptions are not included in this book.

## *Getting started with R and RStudio*

---

### 1.1 Why R?

R is not a programming language like C or Java. It was not created by software engineers for software development. Instead, it was developed by statisticians as an interactive environment for data analysis. You can read the full history in the paper *A Brief History of S*<sup>1</sup>. The interactivity is an indispensable feature in data science because, as you will soon learn, the ability to quickly explore data is a necessity for success in this field. However, like in other programming languages, you can save your work as scripts that can be easily executed at any moment. These scripts serve as a record of the analysis you performed, a key feature that facilitates reproducible work. If you are an expert programmer, you should not expect R to follow the conventions you are used to since you will be disappointed. If you are patient, you will come to appreciate the unequal power of R when it comes to data analysis and, specifically, data visualization.

Other attractive features of R are:

1. R is free and open source<sup>2</sup>.
2. It runs on all major platforms: Windows, Mac Os, UNIX/Linux.
3. Scripts and data objects can be shared seamlessly across platforms.
4. There is a large, growing, and active community of R users and, as a result, there are numerous resources for learning and asking questions<sup>3 4 5</sup>.
5. It is easy for others to contribute add-ons which enables developers to share software implementations of new data science methodologies. This gives R users early access to the latest methods and to tools which are developed for a wide variety of disciplines, including ecology, molecular biology, social sciences, and geography, just to name a few examples.

---

### 1.2 The R console

Interactive data analysis usually occurs on the *R console* that executes commands as you type them. There are several ways to gain access to an R console. One way is to simply start R on your computer. The console looks something like this:

---

<sup>1</sup><https://pdfs.semanticscholar.org/9b48/46f192aa37ca122cfabb1ed1b59866d8bfda.pdf>

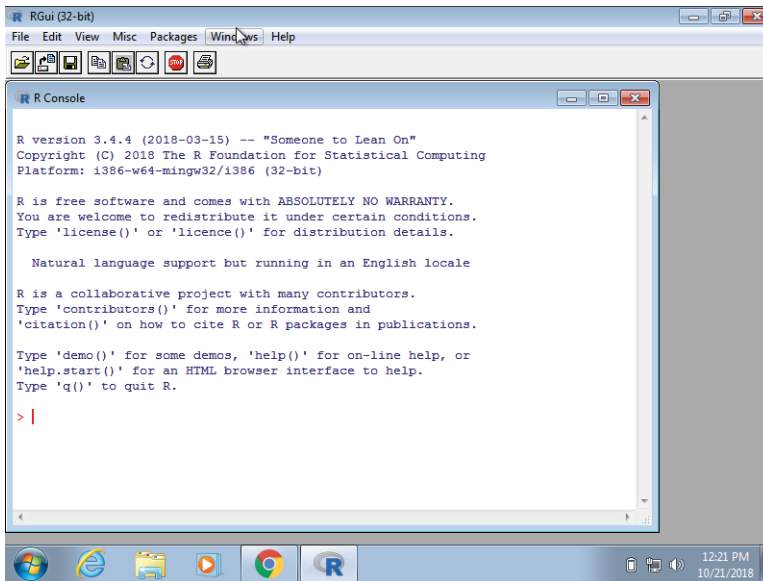
<sup>2</sup><https://opensource.org/history>

<sup>3</sup><https://stats.stackexchange.com/questions/138/free-resources-for-learning-r>

<sup>4</sup><https://www.r-project.org/help.html>

<sup>5</sup><https://stackoverflow.com/documentation/r/topics>





As a quick example, try using the console to calculate a 15% tip on a meal that cost \$19.71:

```
0.15 * 19.71
#> [1] 2.96
```

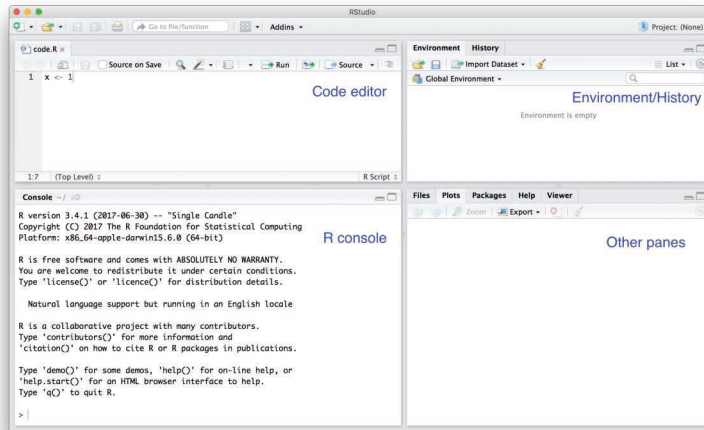
Note that in this book, grey boxes are used to show R code typed into the R console. The symbol `#>` is used to denote what the R console outputs.

## 1.3 Scripts

One of the great advantages of R over point-and-click analysis software is that you can save your work as scripts. You can edit and save these scripts using a text editor. The material in this book was developed using the interactive *integrated development environment* (IDE) RStudio<sup>6</sup>. RStudio includes an editor with many R specific features, a console to execute your code, and other useful panes, including one to show figures.

---

<sup>6</sup><https://www.rstudio.com/>



Most web-based R consoles also provide a pane to edit scripts, but not all permit you to save the scripts for later use.

All the R scripts used to generate this book can be found on GitHub<sup>7</sup>.

---

## 1.4 RStudio

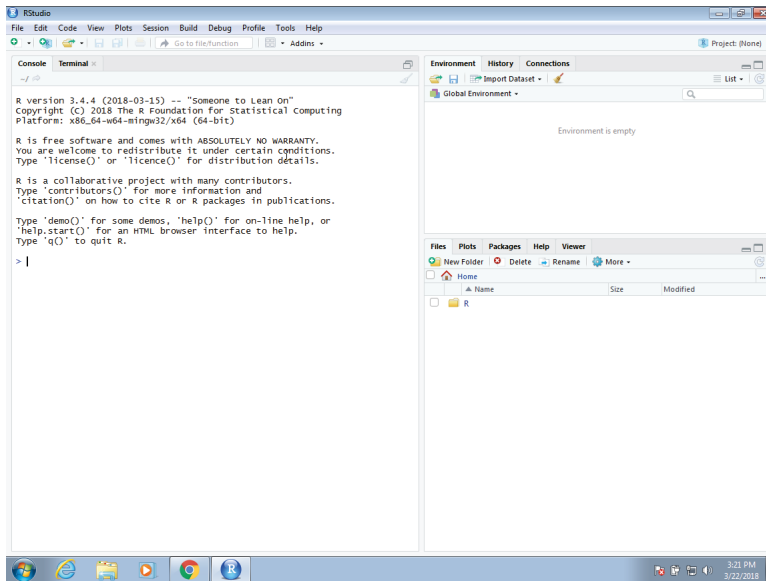
RStudio will be our launching pad for data science projects. It not only provides an editor for us to create and edit our scripts but also provides many other useful tools. In this section, we go over some of the basics.

### 1.4.1 The panes

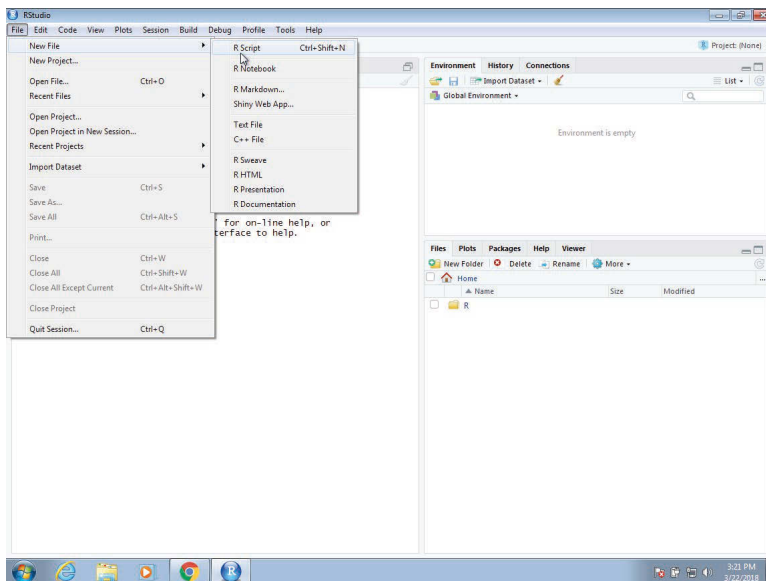
When you start RStudio for the first time, you will see three panes. The left pane shows the R console. On the right, the top pane includes tabs such as *Environment* and *History*, while the bottom pane shows five tabs: *File*, *Plots*, *Packages*, *Help*, and *Viewer* (these tabs may change in new versions). You can click on each tab to move across the different features.

---

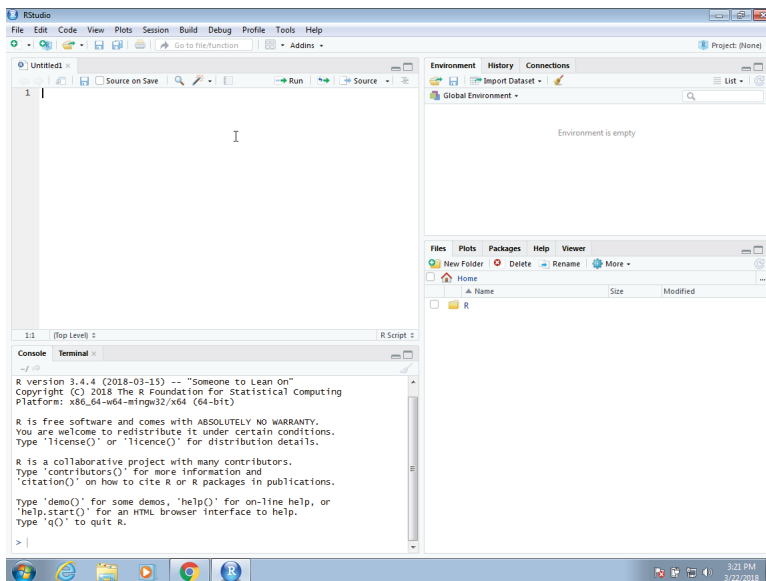
<sup>7</sup><https://github.com/rafalab/dsbook>



To start a new script, you can click on File, the New File, then R Script.



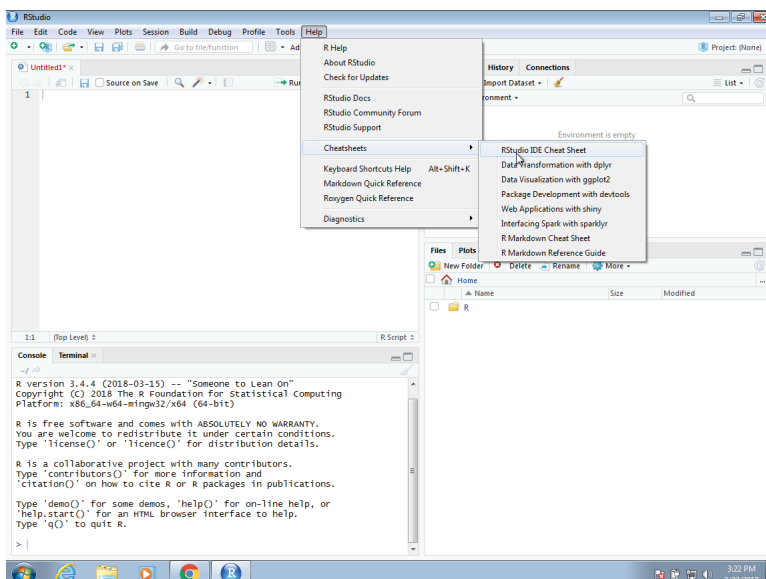
This starts a new pane on the left and it is here where you can start writing your script.



### 1.4.2 Key bindings

Many tasks we perform with the mouse can be achieved with a combination of key strokes instead. These keyboard versions for performing tasks are referred to as *key bindings*. For example, we just showed how to use the mouse to start a new script, but you can also use a key binding: Ctrl+Shift+N on Windows and command+shift+N on the Mac.

Although in this tutorial we often show how to use the mouse, **we highly recommend that you memorize key bindings for the operations you use most**. RStudio provides a useful cheat sheet with the most widely used commands. You can get it from RStudio directly:



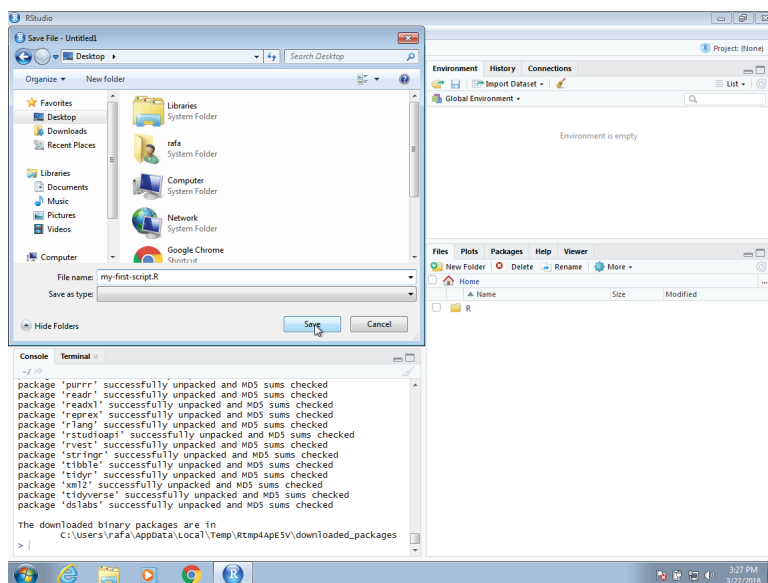
You might want to keep this handy so you can look up key-bindings when you find yourself performing repetitive point-and-clicking.

### 1.4.3 Running commands while editing scripts

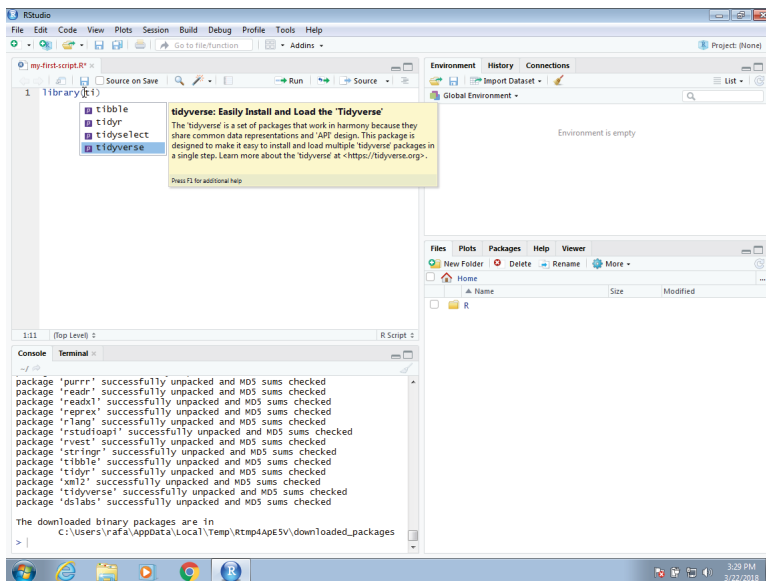
There are many editors specifically made for coding. These are useful because color and indentation are automatically added to make code more readable. RStudio is one of these editors, and it was specifically developed for R. One of the main advantages provided by RStudio over other editors is that we can test our code easily as we edit our scripts. Below we show an example.

Let's start by opening a new script as we did before. A next step is to give the script a name. We can do this through the editor by saving the current new unnamed script. To do this, click on the save icon or use the key binding `Ctrl+S` on Windows and `command+S` on the Mac.

When you ask for the document to be saved for the first time, RStudio will prompt you for a name. A good convention is to use a descriptive name, with lower case letters, no spaces, only hyphens to separate words, and then followed by the suffix `.R`. We will call this script *my-first-script.R*.



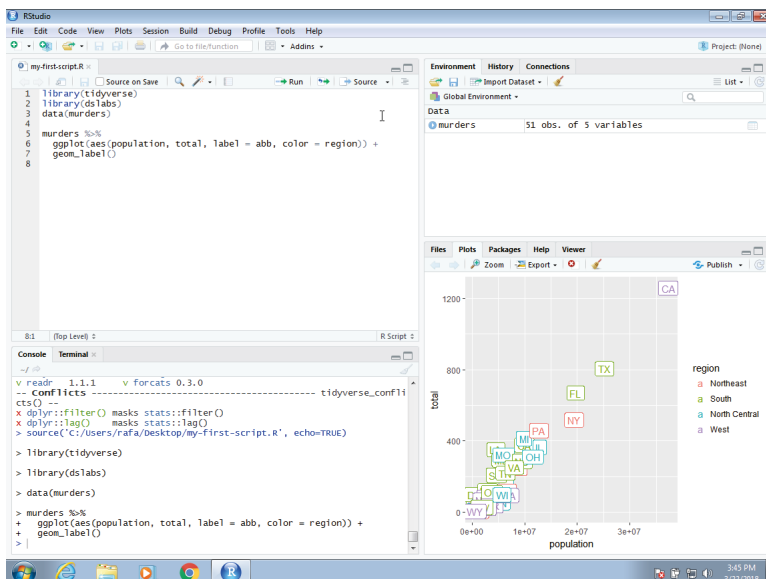
Now we are ready to start editing our first script. The first lines of code in an R script are dedicated to loading the libraries we will use. Another useful RStudio feature is that once we type `library()` it starts auto-completing with libraries that we have installed. Note what happens when we type `library(ti)`:



Another feature you may have noticed is that when you type `library(` the second parenthesis is automatically added. This will help you avoid one of the most common errors in coding: forgetting to close a parenthesis.

Now we can continue to write code. As an example, we will make a graph showing murder totals versus population totals by state. Once you are done writing the code needed to make this plot, you can try it out by *executing* the code. To do this, click on the *Run* button on the upper right side of the editing pane. You can also use the key binding: `Ctrl+Shift+Enter` on Windows or `command+shift+return` on the Mac.

Once you run the code, you will see it appear in the R console and, in this case, the generated plot appears in the plots console. Note that the plot console has a useful interface that permits you to click back and forward across different plots, zoom in to the plot, or save the plots as files.



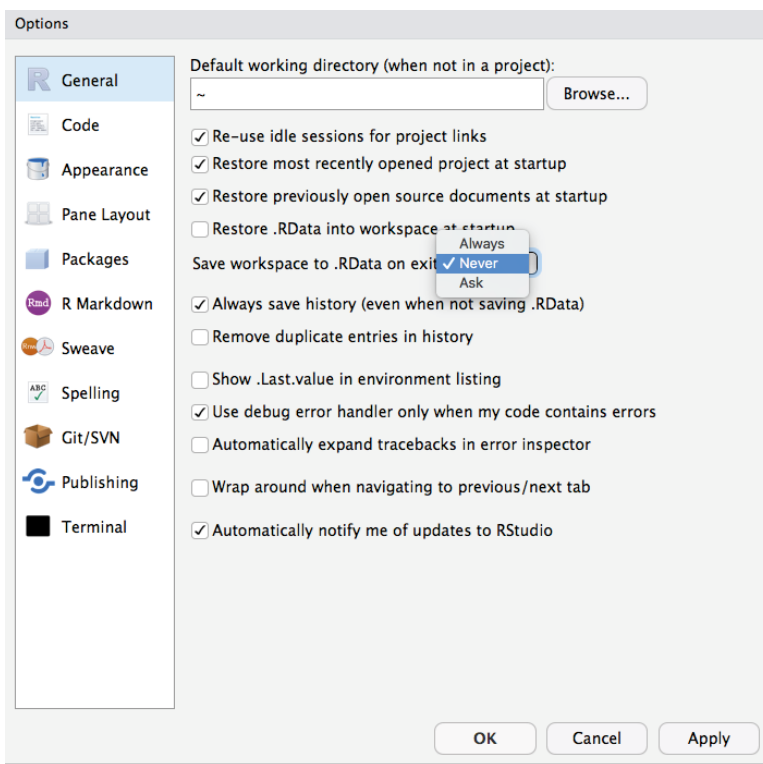
To run one line at a time instead of the entire script, you can use Control-Enter on Windows and command-return on the Mac.

### 1.4.4 Changing global options

You can change the look and functionality of RStudio quite a bit.

To change the global options you click on *Tools* then *Global Options...*

As an example we show how to make a change that we **highly recommend**. This is to change the *Save workspace to .RData on exit* to *Never* and uncheck the *Restore .RData into workspace at start*. By default, when you exit R saves all the objects you have created into a file called .RData. This is done so that when you restart the session in the same folder, it will load these objects. We find that this causes confusion especially when we share code with colleagues and assume they have this .RData file. To change these options, make your *General* settings look like this:



## 1.5 Installing R packages

The functionality provided by a fresh install of R is only a small fraction of what is possible. In fact, we refer to what you get after your first install as *base R*. The extra functionality comes from add-ons available from developers. There are currently hundreds of these available from

CRAN and many others shared via other repositories such as GitHub. However, because not everybody needs all available functionality, R instead makes different components available via *packages*. R makes it very easy to install packages from within R. For example, to install the **dslabs** package, which we use to share datasets and code related to this book, you would type:

```
install.packages("dslabs")
```

In RStudio, you can navigate to the *Tools* tab and select install packages. We can then load the package into our R sessions using the **library** function:

```
library(dslabs)
```

As you go through this book, you will see that we load packages without installing them. This is because once you install a package, it remains installed and only needs to be loaded with **library**. The package remains loaded until we quit the R session. If you try to load a package and get an error, it probably means you need to install it first.

We can install more than one package at once by feeding a character vector to this function:

```
install.packages(c("tidyverse", "dslabs"))
```

Note that installing **tidyverse** actually installs several packages. This commonly occurs when a package has *dependencies*, or uses functions from other packages. When you load a package using **library**, you also load its dependencies.

Once packages are installed, you can load them into R and you do not need to install them again, unless you install a fresh version of R. Remember packages are installed in R not RStudio.

It is helpful to keep a list of all the packages you need for your work in a script because if you need to perform a fresh install of R, you can re-install all your packages by simply running a script.

You can see all the packages you have installed using the following function:

```
installed.packages()
```





# Taylor & Francis

Taylor & Francis Group

<http://taylorandfrancis.com>

# Part I

## R



# Taylor & Francis

Taylor & Francis Group

<http://taylorandfrancis.com>

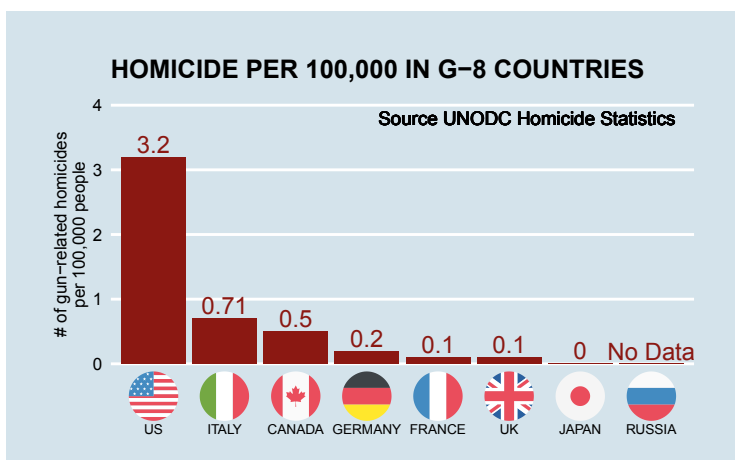
# 2

## *R* basics

In this book, we will be using the R software environment for all our analysis. You will learn R and data analysis techniques simultaneously. To follow along you will therefore need access to R. We also recommend the use of an *integrated development environment* (IDE), such as RStudio, to save your work. Note that it is common for a course or workshop to offer access to an R environment and an IDE through your web browser, as done by RStudio cloud<sup>1</sup>. If you have access to such a resource, you don't need to install R and RStudio. However, if you intend on becoming an advanced data analyst, we highly recommend installing these tools on your computer<sup>2</sup>. Both R and RStudio are free and available online.

### 2.1 Case study: US Gun Murders

Imagine you live in Europe and are offered a job in a US company with many locations across all states. It is a great job, but news with headlines such as **US Gun Homicide Rate Higher Than Other Developed Countries**<sup>3</sup> have you worried. Charts like this may concern you even more:

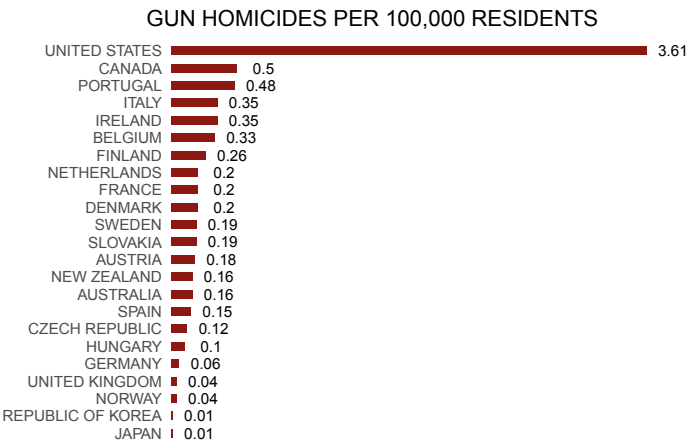


<sup>1</sup><https://rstudio.cloud>

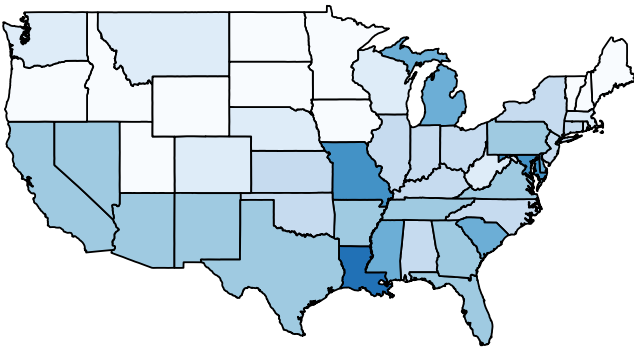
<sup>2</sup><https://rafalab.github.io/dsbook/installing-r-rstudio.html>

<sup>3</sup><http://abcnews.go.com/blogs/headlines/2012/12/us-gun-ownership-homicide-rate-higher-than-other-developed-countries/>

Or even worse, this version from [everytown.org](http://everytown.org):



But then you remember that the US is a large and diverse country with 50 very different states as well as the District of Columbia (DC).



California, for example, has a larger population than Canada, and 20 US states have populations larger than that of Norway. In some respects, the variability across states in the US is akin to the variability across countries in Europe. Furthermore, although not included in the charts above, the murder rates in Lithuania, Ukraine, and Russia are higher than 4 per 100,000. So perhaps the news reports that worried you are too superficial. You have options of where to live and want to determine the safety of each particular state. We will gain some insights by examining data related to gun homicides in the US during 2010 using R.

Before we get started with our example, we need to cover logistics as well as some of the very basic building blocks that are required to gain more advanced R skills. Be aware that the usefulness of some of these building blocks may not be immediately obvious, but later in the book you will appreciate having mastered these skills.

---

## 2.2 The very basics

Before we get started with the motivating dataset, we need to cover the very basics of R.

### 2.2.1 Objects

Suppose a high school student asks us for help solving several quadratic equations of the form  $ax^2 + bx + c = 0$ . The quadratic formula gives us the solutions:

$$\frac{-b - \sqrt{b^2 - 4ac}}{2a} \quad \text{and} \quad \frac{-b + \sqrt{b^2 - 4ac}}{2a}$$

which of course change depending on the values of  $a$ ,  $b$ , and  $c$ . One advantage of programming languages is that we can define variables and write expressions with these variables, similar to how we do so in math, but obtain a numeric solution. We will write out general code for the quadratic equation below, but if we are asked to solve  $x^2 + x - 1 = 0$ , then we define:

```
a <- 1
b <- 1
c <- -1
```

which stores the values for later use. We use `<-` to assign values to the variables.

We can also assign values using `=` instead of `<-`, but we recommend against using `=` to avoid confusion.

Copy and paste the code above into your console to define the three variables. Note that R does not print anything when we make this assignment. This means the objects were defined successfully. Had you made a mistake, you would have received an error message.

To see the value stored in a variable, we simply ask R to evaluate `a` and it shows the stored value:

```
a
#> [1] 1
```

A more explicit way to ask R to show us the value stored in `a` is using `print` like this:

```
print(a)
#> [1] 1
```

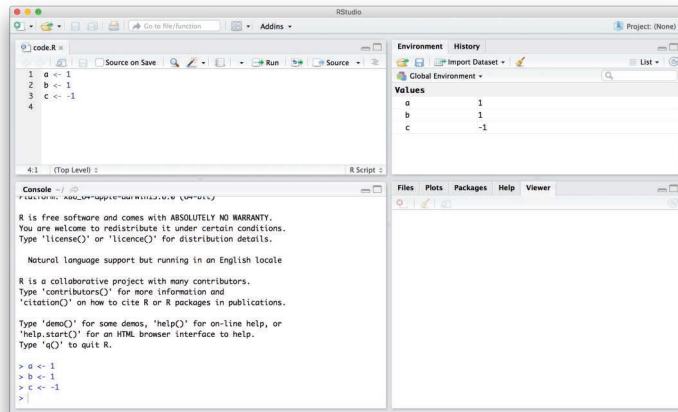
We use the term *object* to describe stuff that is stored in R. Variables are examples, but objects can also be more complicated entities such as functions, which are described later.

### 2.2.2 The workspace

As we define objects in the console, we are actually changing the *workspace*. You can see all the variables saved in your workspace by typing:

```
ls()
#> [1] "a"      "b"      "c"      "dat"    "img_path" "murders"
```

In RStudio, the *Environment* tab shows the values:



We should see `a`, `b`, and `c`. If you try to recover the value of a variable that is not in your workspace, you receive an error. For example, if you type `x` you will receive the following message: `Error: object 'x' not found`.

Now since these values are saved in variables, to obtain a solution to our equation, we use the quadratic formula:

```
(-b + sqrt(b^2 - 4*a*c) ) / ( 2*a )
#> [1] 0.618
(-b - sqrt(b^2 - 4*a*c) ) / ( 2*a )
#> [1] -1.62
```

### 2.2.3 Functions

Once you define variables, the data analysis process can usually be described as a series of *functions* applied to the data. R includes several predefined functions and most of the analysis pipelines we construct make extensive use of these.

We already used the `install.packages`, `library`, and `ls` functions. We also used the function `sqrt` to solve the quadratic equation above. There are many more prebuilt functions and even more can be added through packages. These functions do not appear in the workspace because you did not define them, but they are available for immediate use.

In general, we need to use parentheses to evaluate a function. If you type `ls`, the function is not evaluated and instead R shows you the code that defines the function. If you type `ls()` the function is evaluated and, as seen above, we see objects in the workspace.

Unlike `ls`, most functions require one or more *arguments*. Below is an example of how we

assign an object to the argument of the function `log`. Remember that we earlier defined `a` to be 1:

```
log(8)
#> [1] 2.08
log(a)
#> [1] 0
```

You can find out what the function expects and what it does by reviewing the very useful manuals included in R. You can get help by using the `help` function like this:

```
help("log")
```

For most functions, we can also use this shorthand:

```
?log
```

The help page will show you what arguments the function is expecting. For example, `log` needs `x` and `base` to run. However, some arguments are required and others are optional. You can determine which arguments are optional by noting in the help document that a default value is assigned with `=`. Defining these is optional. For example, the base of the function `log` defaults to `base = exp(1)` making `log` the natural log by default.

If you want a quick look at the arguments without opening the help system, you can type:

```
args(log)
#> function (x, base = exp(1))
#> NULL
```

You can change the default values by simply assigning another object:

```
log(8, base = 2)
#> [1] 3
```

Note that we have not been specifying the argument `x` as such:

```
log(x = 8, base = 2)
#> [1] 3
```

The above code works, but we can save ourselves some typing: if no argument name is used, R assumes you are entering arguments in the order shown in the help file or by `args`. So by not using the names, it assumes the arguments are `x` followed by `base`:

```
log(8,2)
#> [1] 3
```

If using the arguments' names, then we can include them in whatever order we want:

```
log(base = 2, x = 8)
#> [1] 3
```



To specify arguments, we must use `=`, and cannot use `<-`.

There are some exceptions to the rule that functions need the parentheses to be evaluated. Among these, the most commonly used are the arithmetic and relational operators. For example:

```
2 ^ 3
#> [1] 8
```

You can see the arithmetic operators by typing:

```
help("+")
```

or

```
? "+"
```

and the relational operators by typing:

```
help(">")
```

or

```
? ">"
```

### 2.2.4 Other prebuilt objects

There are several datasets that are included for users to practice and test out functions. You can see all the available datasets by typing:

```
data()
```

This shows you the object name for these datasets. These datasets are objects that can be used by simply typing the name. For example, if you type:

```
co2
```

R will show you Mauna Loa atmospheric CO2 concentration data.

Other prebuilt objects are mathematical quantities, such as the constant  $\pi$  and  $\infty$ :

```
pi
#> [1] 3.14
Inf+1
#> [1] Inf
```

### 2.2.5 Variable names

We have used the letters `a`, `b`, and `c` as variable names, but variable names can be almost anything. Some basic rules in R are that variable names have to start with a letter, can't

contain spaces, and should not be variables that are predefined in R. For example, don't name one of your variables `install.packages` by typing something like `install.packages <- 2`.

A nice convention to follow is to use meaningful words that describe what is stored, use only lower case, and use underscores as a substitute for spaces. For the quadratic equations, we could use something like this:

```
solution_1 <- (-b + sqrt(b^2 - 4*a*c)) / (2*a)
solution_2 <- (-b - sqrt(b^2 - 4*a*c)) / (2*a)
```

For more advice, we highly recommend studying Hadley Wickham's style guide<sup>4</sup>.

### 2.2.6 Saving your workspace

Values remain in the workspace until you end your session or erase them with the function `rm`. But workspaces also can be saved for later use. In fact, when you quit R, the program asks you if you want to save your workspace. If you do save it, the next time you start R, the program will restore the workspace.

We actually recommend against saving the workspace this way because, as you start working on different projects, it will become harder to keep track of what is saved. Instead, we recommend you assign the workspace a specific name. You can do this by using the function `save` or `save.image`. To load, use the function `load`. When saving a workspace, we recommend the suffix `rda` or `RData`. In RStudio, you can also do this by navigating to the *Session* tab and choosing *Save Workspace as*. You can later load it using the *Load Workspace* options in the same tab. You can read the help pages on `save`, `save.image`, and `load` to learn more.

### 2.2.7 Motivating scripts

To solve another equation such as  $3x^2 + 2x - 1$ , we can copy and paste the code above and then redefine the variables and recompute the solution:

```
a <- 3
b <- 2
c <- -1
(-b + sqrt(b^2 - 4*a*c)) / (2*a)
(-b - sqrt(b^2 - 4*a*c)) / (2*a)
```

By creating and saving a script with the code above, we would not need to retype everything each time and, instead, simply change the variable names. Try writing the script above into an editor and notice how easy it is to change the variables and receive an answer.

---

<sup>4</sup><http://adv-r.had.co.nz/Style.html>

## 2.2.8 Commenting your code

If a line of R code starts with the symbol `#`, it is not evaluated. We can use this to write reminders of why we wrote particular code. For example, in the script above we could add:

```
## Code to compute solution to quadratic equation of the form ax^2 + bx + c
## define the variables
a <- 3
b <- 2
c <- -1

## now compute the solution
(-b + sqrt(b^2 - 4*a*c)) / (2*a)
(-b - sqrt(b^2 - 4*a*c)) / (2*a)
```

---

## 2.3 Exercises

1. What is the sum of the first 100 positive integers? The formula for the sum of integers 1 through  $n$  is  $n(n+1)/2$ . Define  $n = 100$  and then use R to compute the sum of 1 through 100 using the formula. What is the sum?
2. Now use the same formula to compute the sum of the integers from 1 through 1,000.
3. Look at the result of typing the following code into R:

```
n <- 1000
x <- seq(1, n)
sum(x)
```

Based on the result, what do you think the functions `seq` and `sum` do? You can use `help`.

- a. `sum` creates a list of numbers and `seq` adds them up.
  - b. `seq` creates a list of numbers and `sum` adds them up.
  - c. `seq` creates a random list and `sum` computes the sum of 1 through 1,000.
  - d. `sum` always returns the same number.
4. In math and programming, we say that we evaluate a function when we replace the argument with a given number. So if we type `sqrt(4)`, we evaluate the `sqrt` function. In R, you can evaluate a function inside another function. The evaluations happen from the inside out. Use one line of code to compute the log, in base 10, of the square root of 100.
  5. Which of the following will always return the numeric value stored in `x`? You can try out examples and use the help system if you want.
    - a. `log(10^x)`
    - b. `log10(x^10)`
    - c. `log(exp(x))`
    - d. `exp(log(x, base = 2))`

---

## 2.4 Data types

Variables in R can be of different types. For example, we need to distinguish numbers from character strings and tables from simple lists of numbers. The function `class` helps us determine what type of object we have:

```
a <- 2
class(a)
#> [1] "numeric"
```

To work efficiently in R, it is important to learn the different types of variables and what we can do with these.

### 2.4.1 Data frames

Up to now, the variables we have defined are just one number. This is not very useful for storing data. The most common way of storing a dataset in R is in a *data frame*. Conceptually, we can think of a data frame as a table with rows representing observations and the different variables reported for each observation defining the columns. Data frames are particularly useful for datasets because we can combine different data types into one object.

A large proportion of data analysis challenges start with data stored in a data frame. For example, we stored the data for our motivating example in a data frame. You can access this dataset by loading the **dslabs** library and loading the **murders** dataset using the `data` function:

```
library(dslabs)
data(murders)
```

To see that this is in fact a data frame, we type:

```
class(murders)
#> [1] "data.frame"
```

### 2.4.2 Examining an object

The function `str` is useful for finding out more about the structure of an object:

```
str(murders)
#> 'data.frame':   51 obs. of  5 variables:
#> $ state : chr "Alabama" "Alaska" "Arizona" "Arkansas" ...
#> $ abb  : chr "AL" "AK" "AZ" "AR" ...
#> $ region : Factor w/ 4 levels "Northeast","South",...: 2 4 4 2 4 4 1 2 2
#> 2 ...
#> $ population: num 4779736 710231 6392017 2915918 37253956 ...
#> $ total : num 135 19 232 93 1257 ...
```

This tells us much more about the object. We see that the table has 51 rows (50 states plus DC) and five variables. We can show the first six lines using the function `head`:

```
head(murders)
#>      state abb region population total
#> 1  Alabama AL  South    4779736    135
#> 2  Alaska  AK   West     710231     19
#> 3  Arizona AZ   West    6392017    232
#> 4  Arkansas AR  South    2915918     93
#> 5 California CA   West    37253956  1257
#> 6  Colorado CO   West    5029196     65
```

In this dataset, each state is considered an observation and five variables are reported for each state.

Before we go any further in answering our original question about different states, let's learn more about the components of this object.

### 2.4.3 The accessor: `$`

For our analysis, we will need to access the different variables represented by columns included in this data frame. To do this, we use the accessor operator `$` in the following way:

```
murders$population
#> [1] 4779736 710231 6392017 2915918 37253956 5029196 3574097
#> [8] 897934 601723 19687653 9920000 1360301 1567582 12830632
#> [15] 6483802 3046355 2853118 4339367 4533372 1328361 5773552
#> [22] 6547629 9883640 5303925 2967297 5988927 989415 1826341
#> [29] 2700551 1316470 8791894 2059179 19378102 9535483 672591
#> [36] 11536504 3751351 3831074 12702379 1052567 4625364 814180
#> [43] 6346105 25145561 2763885 625741 8001024 6724540 1852994
#> [50] 5686986 563626
```

But how did we know to use `population`? Previously, by applying the function `str` to the object `murders`, we revealed the names for each of the five variables stored in this table. We can quickly access the variable names using:

```
names(murders)
#> [1] "state"      "abb"        "region"     "population" "total"
```

It is important to know that the order of the entries in `murders$population` preserves the order of the rows in our data table. This will later permit us to manipulate one variable based on the results of another. For example, we will be able to order the state names by the number of murders.

**Tip:** R comes with a very nice auto-complete functionality that saves us the trouble of typing out all the names. Try typing `murders$p` then hitting the `tab` key on your keyboard. This functionality and many other useful auto-complete features are available when working in RStudio.

### 2.4.4 Vectors: numerics, characters, and logical

The object `murders$population` is not one number but several. We call these types of objects *vectors*. A single number is technically a vector of length 1, but in general we use the term vectors to refer to objects with several entries. The function `length` tells you how many entries are in the vector:

```
pop <- murders$population
length(pop)
#> [1] 51
```

This particular vector is *numeric* since population sizes are numbers:

```
class(pop)
#> [1] "numeric"
```

In a numeric vector, every entry must be a number.

To store character strings, vectors can also be of class *character*. For example, the state names are characters:

```
class(murders$state)
#> [1] "character"
```

As with numeric vectors, all entries in a character vector need to be a character.

Another important type of vectors are *logical vectors*. These must be either `TRUE` or `FALSE`.

```
z <- 3 == 2
z
#> [1] FALSE
class(z)
#> [1] "logical"
```

Here the `==` is a relational operator asking if 3 is equal to 2. In R, if you just use one `=`, you actually assign a variable, but if you use two `==` you test for equality.

You can see the other *relational operators* by typing:

```
?Comparison
```

In future sections, you will see how useful relational operators can be.

We discuss more important features of vectors after the next set of exercises.

**Advanced:** Mathematically, the values in `pop` are integers and there is an integer class in R. However, by default, numbers are assigned class `numeric` even when they are round integers. For example, `class(1)` returns `numeric`. You can turn them into class `integer` with the `as.integer()` function or by adding an `L` like this: `1L`. Note the class by typing: `class(1L)`

### 2.4.5 Factors

In the `murders` dataset, we might expect the `region` to also be a character vector. However, it is not:

```
class(murders$region)
#> [1] "factor"
```

It is a *factor*. Factors are useful for storing categorical data. We can see that there are only 4 regions by using the `levels` function:

```
levels(murders$region)
#> [1] "Northeast"      "South"           "North Central"  "West"
```

In the background, R stores these *levels* as integers and keeps a map to keep track of the labels. This is more memory efficient than storing all the characters.

Note that the levels have an order that is different from the order of appearance in the factor object. The default is for the levels to follow alphabetical order. However, often we want the levels to follow a different order. We will see several examples of this in the Data Visualization part of the book. The function `reorder` lets us change the order of the levels of a factor variable based on a summary computed on a numeric vector. We will demonstrate this with a simple example.

Suppose we want the levels of the `region` by the total number of murders rather than alphabetical order. If there are values associated with each level, we can use the `reorder` and specify a data summary to determine the order. The following code takes the sum of the total murders in each region, and reorders the factor following these sums.

```
region <- murders$region
value <- murders$total
region <- reorder(region, value, FUN = sum)
levels(region)
#> [1] "Northeast"      "North Central"  "West"           "South"
```

The new order is in agreement with the fact that the Northeast has the least murders and the South has the most.

**Warning:** Factors can be a source of confusion since sometimes they behave like characters and sometimes they do not. As a result, confusing factors and characters are a common source of bugs.

### 2.4.6 Lists

Data frames are a special case of *lists*. We will cover lists in more detail later, but know that they are useful because you can store any combination of different types. Below is an example of a list we created for you:

```
record
#> $name
#> [1] "John Doe"
#>
#> $student_id
#> [1] 1234
#>
#> $grades
#> [1] 95 82 91 97 93
#>
#> $final_grade
#> [1] "A"
class(record)
#> [1] "list"
```

As with data frames, you can extract the components of a list with the accessor `$`. In fact, data frames are a type of list.

```
record$student_id
#> [1] 1234
```

We can also use double square brackets (`[[`) like this:

```
record[["student_id"]]
#> [1] 1234
```

You should get used to the fact that in R, there are often several ways to do the same thing, such as accessing entries.

You might also encounter lists without variable names.

```
record2
#> [[1]]
#> [1] "John Doe"
#>
#> [[2]]
#> [1] 1234
```

If a list does not have names, you cannot extract the elements with `$`, but you can still use the brackets method and instead of providing the variable name, you provide the list index, like this:

```
record2[[1]]
#> [1] "John Doe"
```

We won't be using lists until later, but you might encounter one in your own exploration of R. For this reason, we show you some basics here.



### 2.4.7 Matrices

Matrices are another type of object that are common in R. Matrices are similar to data frames in that they are two-dimensional: they have rows and columns. However, like numeric, character and logical vectors, entries in matrices have to be all the same type. For this reason data frames are much more useful for storing data, since we can have characters, factors, and numbers in them.

Yet matrices have a major advantage over data frames: we can perform matrix algebra operations, a powerful type of mathematical technique. We do not describe these operations in this book, but much of what happens in the background when you perform a data analysis involves matrices. We cover matrices in more detail in Chapter 33.1 but describe them briefly here since some of the functions we will learn return matrices.

We can define a matrix using the `matrix` function. We need to specify the number of rows and columns.

```
mat <- matrix(1:12, 4, 3)
mat
#>      [,1] [,2] [,3]
#> [1,]    1    5    9
#> [2,]    2    6   10
#> [3,]    3    7   11
#> [4,]    4    8   12
```

You can access specific entries in a matrix using square brackets (`[]`). If you want the second row, third column, you use:

```
mat[2, 3]
#> [1] 10
```

If you want the entire second row, you leave the column spot empty:

```
mat[2, ]
#> [1]  2  6 10
```

Notice that this returns a vector, not a matrix.

Similarly, if you want the entire third column, you leave the row spot empty:

```
mat[, 3]
#> [1]  9 10 11 12
```

This is also a vector, not a matrix.

You can access more than one column or more than one row if you like. This will give you a new matrix.

```
mat[, 2:3]
#>      [,1] [,2]
#> [1,]    5    9
#> [2,]    6   10
```

```
#> [3,]    7    11
#> [4,]    8    12
```

You can subset both rows and columns:

```
mat[1:2, 2:3]
#>      [,1] [,2]
#> [1,]    5    9
#> [2,]    6   10
```

We can convert matrices into data frames using the function `as.data.frame`:

```
as.data.frame(mat)
#>   V1 V2 V3
#> 1  1  5  9
#> 2  2  6 10
#> 3  3  7 11
#> 4  4  8 12
```

You can also use single square brackets (`[`) to access rows and columns of a data frame:

```
data("murders")
murders[25, 1]
#> [1] "Mississippi"
murders[2:3, ]
#>   state abb region population total
#> 2 Alaska AK  West    710231     19
#> 3 Arizona AZ  West    6392017    232
```

---

## 2.5 Exercises

1. Load the US murders dataset.

```
library(dslabs)
data(murders)
```

Use the function `str` to examine the structure of the `murders` object. Which of the following best describes the variables represented in this data frame?

- The 51 states.
- The murder rates for all 50 states and DC.
- The state name, the abbreviation of the state name, the state's region, and the state's population and total number of murders for 2010.
- `str` shows no relevant information.

2. What are the column names used by the data frame for these five variables?
3. Use the accessor `$` to extract the state abbreviations and assign them to the object `a`. What is the class of this object?
4. Now use the square brackets to extract the state abbreviations and assign them to the object `b`. Use the `identical` function to determine if `a` and `b` are the same.
5. We saw that the `region` column stores a factor. You can corroborate this by typing:

```
class(murders$region)
```

With one line of code, use the function `levels` and `length` to determine the number of regions defined by this dataset.

6. The function `table` takes a vector and returns the frequency of each element. You can quickly see how many states are in each region by applying this function. Use this function in one line of code to create a table of states per region.

## 2.6 Vectors

In R, the most basic objects available to store data are *vectors*. As we have seen, complex datasets can usually be broken down into components that are vectors. For example, in a data frame, each column is a vector. Here we learn more about this important class.

### 2.6.1 Creating vectors

We can create vectors using the function `c`, which stands for *concatenate*. We use `c` to concatenate entries in the following way:

```
codes <- c(380, 124, 818)
codes
#> [1] 380 124 818
```

We can also create character vectors. We use the quotes to denote that the entries are characters rather than variable names.

```
country <- c("italy", "canada", "egypt")
```

In R you can also use single quotes:

```
country <- c('italy', 'canada', 'egypt')
```

But be careful not to confuse the single quote `'` with the *back quote* ```.

By now you should know that if you type:

```
country <- c(italy, canada, egypt)
```

you receive an error because the variables `italy`, `canada`, and `egypt` are not defined. If we do not use the quotes, R looks for variables with those names and returns an error.

### 2.6.2 Names

Sometimes it is useful to name the entries of a vector. For example, when defining a vector of country codes, we can use the names to connect the two:

```
codes <- c(italy = 380, canada = 124, egypt = 818)
codes
#>   italy  canada  egypt
#>   380    124    818
```

The object `codes` continues to be a numeric vector:

```
class(codes)
#> [1] "numeric"
```

but with names:

```
names(codes)
#> [1] "italy" "canada" "egypt"
```

If the use of strings without quotes looks confusing, know that you can use the quotes as well:

```
codes <- c("italy" = 380, "canada" = 124, "egypt" = 818)
codes
#>   italy  canada  egypt
#>   380    124    818
```

There is no difference between this function call and the previous one. This is one of the many ways in which R is quirky compared to other languages.

We can also assign names using the `names` functions:

```
codes <- c(380, 124, 818)
country <- c("italy", "canada", "egypt")
names(codes) <- country
codes
#>   italy  canada  egypt
#>   380    124    818
```

### 2.6.3 Sequences

Another useful function for creating vectors generates sequences:

```
seq(1, 10)
#> [1] 1 2 3 4 5 6 7 8 9 10
```

The first argument defines the start, and the second defines the end which is included. The default is to go up in increments of 1, but a third argument lets us tell it how much to jump by:

```
seq(1, 10, 2)
#> [1] 1 3 5 7 9
```

If we want consecutive integers, we can use the following shorthand:

```
1:10
#> [1] 1 2 3 4 5 6 7 8 9 10
```

When we use these functions, R produces integers, not numerics, because they are typically used to index something:

```
class(1:10)
#> [1] "integer"
```

However, if we create a sequence including non-integers, the class changes:

```
class(seq(1, 10, 0.5))
#> [1] "numeric"
```

## 2.6.4 Subsetting

We use square brackets to access specific elements of a vector. For the vector `codes` we defined above, we can access the second element using:

```
codes[2]
#> canada
#> 124
```

You can get more than one entry by using a multi-entry vector as an index:

```
codes[c(1,3)]
#> italy egypt
#> 380 818
```

The sequences defined above are particularly useful if we want to access, say, the first two elements:

```
codes[1:2]
#> italy canada
#> 380 124
```

If the elements have names, we can also access the entries using these names. Below are two examples.

```
codes["canada"]
#> canada
```

```
#>      124
codes[c("egypt","italy")]
#> egypt  italy
#>      818    380
```

## 2.7 Coercion

In general, *coercion* is an attempt by R to be flexible with data types. When an entry does not match the expected, some of the prebuilt R functions try to guess what was meant before throwing an error. This can also lead to confusion. Failing to understand *coercion* can drive programmers crazy when attempting to code in R since it behaves quite differently from most other languages in this regard. Let's learn about it with some examples.

We said that vectors must be all of the same type. So if we try to combine, say, numbers and characters, you might expect an error:

```
x <- c(1, "canada", 3)
```

But we don't get one, not even a warning! What happened? Look at `x` and its class:

```
x
#> [1] "1"      "canada" "3"
class(x)
#> [1] "character"
```

R *coerced* the data into characters. It guessed that because you put a character string in the vector, you meant the 1 and 3 to actually be character strings "1" and "3". The fact that not even a warning is issued is an example of how coercion can cause many unnoticed errors in R.

R also offers functions to change from one type to another. For example, you can turn numbers into characters with:

```
x <- 1:5
y <- as.character(x)
y
#> [1] "1" "2" "3" "4" "5"
```

You can turn it back with `as.numeric`:

```
as.numeric(y)
#> [1] 1 2 3 4 5
```

This function is actually quite useful since datasets that include numbers as character strings are common.

### 2.7.1 Not availables (NA)

When a function tries to coerce one type to another and encounters an impossible case, it usually gives us a warning and turns the entry into a special value called an **NA** for “not available”. For example:

```
x <- c("1", "b", "3")
as.numeric(x)
#> Warning: NAs introduced by coercion
#> [1] 1 NA 3
```

R does not have any guesses for what number you want when you type **b**, so it does not try.

As a data scientist you will encounter the **NA**s often as they are generally used for missing data, a common problem in real-world datasets.

---

## 2.8 Exercises

1. Use the function `c` to create a vector with the average high temperatures in January for Beijing, Lagos, Paris, Rio de Janeiro, San Juan, and Toronto, which are 35, 88, 42, 84, 81, and 30 degrees Fahrenheit. Call the object `temp`.
2. Now create a vector with the city names and call the object `city`.
3. Use the `names` function and the objects defined in the previous exercises to associate the temperature data with its corresponding city.
4. Use the `[` and `:` operators to access the temperature of the first three cities on the list.
5. Use the `[` operator to access the temperature of Paris and San Juan.
6. Use the `:` operator to create a sequence of numbers 12, 13, 14, ..., 73.
7. Create a vector containing all the positive odd numbers smaller than 100.
8. Create a vector of numbers that starts at 6, does not pass 55, and adds numbers in increments of  $4/7$ : 6,  $6 + 4/7$ ,  $6 + 8/7$ , and so on. How many numbers does the list have? Hint: use `seq` and `length`.
9. What is the class of the following object `a <- seq(1, 10, 0.5)`?
10. What is the class of the following object `a <- seq(1, 10)`?
11. The class of `class(a<-1)` is numeric, not integer. R defaults to numeric and to force an integer, you need to add the letter L. Confirm that the class of `1L` is integer.
12. Define the following vector:

```
x <- c("1", "3", "5")
```

and coerce it to get integers.

## 2.9 Sorting

Now that we have mastered some basic R knowledge, let's try to gain some insights into the safety of different states in the context of gun murders.

### 2.9.1 `sort`

Say we want to rank the states from least to most gun murders. The function `sort` sorts a vector in increasing order. We can therefore see the largest number of gun murders by typing:

```
library(dslabs)
data(murders)
sort(murders$total)
#> [1] 2 4 5 5 7 8 11 12 12 16 19 21 22
#> [14] 27 32 36 38 53 63 65 67 84 93 93 97 97
#> [27] 99 111 116 118 120 135 142 207 219 232 246 250 286
#> [40] 293 310 321 351 364 376 413 457 517 669 805 1257
```

However, this does not give us information about which states have which murder totals. For example, we don't know which state had 1257.

### 2.9.2 `order`

The function `order` is closer to what we want. It takes a vector as input and returns the vector of indexes that sorts the input vector. This may sound confusing so let's look at a simple example. We can create a vector and sort it:

```
x <- c(31, 4, 15, 92, 65)
sort(x)
#> [1] 4 15 31 65 92
```

Rather than sort the input vector, the function `order` returns the index that sorts input vector:

```
index <- order(x)
x[index]
#> [1] 4 15 31 65 92
```

This is the same output as that returned by `sort(x)`. If we look at this index, we see why it works:

```
x
#> [1] 31 4 15 92 65
order(x)
#> [1] 2 3 1 5 4
```



The second entry of `x` is the smallest, so `order(x)` starts with 2. The next smallest is the third entry, so the second entry is 3 and so on.

How does this help us order the states by murders? First, remember that the entries of vectors you access with `$` follow the same order as the rows in the table. For example, these two vectors containing state names and abbreviations, respectively, are matched by their order:

```
murders$state[1:6]
#> [1] "Alabama" "Alaska" "Arizona" "Arkansas" "California"
#> [6] "Colorado"
murders$abb[1:6]
#> [1] "AL" "AK" "AZ" "AR" "CA" "CO"
```

This means we can order the state names by their total murders. We first obtain the index that orders the vectors according to murder totals and then index the state names vector:

```
ind <- order(murders$total)
murders$abb[ind]
#> [1] "VT" "ND" "NH" "WY" "HI" "SD" "ME" "ID" "MT" "RI" "AK" "IA" "UT"
#> [14] "WV" "NE" "OR" "DE" "MN" "KS" "CO" "NM" "NV" "AR" "WA" "CT" "WI"
#> [27] "DC" "OK" "KY" "MA" "MS" "AL" "IN" "SC" "TN" "AZ" "NJ" "VA" "NC"
#> [40] "MD" "OH" "MO" "LA" "IL" "GA" "MI" "PA" "NY" "FL" "TX" "CA"
```

According to the above, California had the most murders.

### 2.9.3 max and which.max

If we are only interested in the entry with the largest value, we can use `max` for the value:

```
max(murders$total)
#> [1] 1257
```

and `which.max` for the index of the largest value:

```
i_max <- which.max(murders$total)
murders$state[i_max]
#> [1] "California"
```

For the minimum, we can use `min` and `which.min` in the same way.

Does this mean California is the most dangerous state? In an upcoming section, we argue that we should be considering rates instead of totals. Before doing that, we introduce one last order-related function: `rank`.

### 2.9.4 rank

Although not as frequently used as `order` and `sort`, the function `rank` is also related to order and can be useful. For any given vector it returns a vector with the rank of the first entry, second entry, etc., of the input vector. Here is a simple example:

```
x <- c(31, 4, 15, 92, 65)
rank(x)
#> [1] 3 1 2 5 4
```

To summarize, let's look at the results of the three functions we have introduced:

original	sort	order	rank
31	4	2	3
4	15	3	1
15	31	1	2
92	65	5	5
65	92	4	4

### 2.9.5 Beware of recycling

Another common source of unnoticed errors in R is the use of *recycling*. We saw that vectors are added elementwise. So if the vectors don't match in length, it is natural to assume that we should get an error. But we don't. Notice what happens:

```
x <- c(1,2,3)
y <- c(10, 20, 30, 40, 50, 60, 70)
x+y
#> Warning in x + y: longer object length is not a multiple of shorter
#> object length
#> [1] 11 22 33 41 52 63 71
```

We do get a warning, but no error. For the output, R has recycled the numbers in **x**. Notice the last digit of numbers in the output.

---

## 2.10 Exercises

For these exercises we will use the US murders dataset. Make sure you load it prior to starting.

```
library(dslabs)
data("murders")
```

1. Use the **\$** operator to access the population size data and store it as the object **pop**. Then use the **sort** function to redefine **pop** so that it is sorted. Finally, use the **[** operator to report the smallest population size.
2. Now instead of the smallest population size, find the index of the entry with the smallest population size. Hint: use **order** instead of **sort**.
3. We can actually perform the same operation as in the previous exercise using the function **which.min**. Write one line of code that does this.
4. Now we know how small the smallest state is and we know which row represents it. Which

state is it? Define a variable `states` to be the state names from the `murders` data frame. Report the name of the state with the smallest population.

5. You can create a data frame using the `data.frame` function. Here is a quick example:

```
temp <- c(35, 88, 42, 84, 81, 30)
city <- c("Beijing", "Lagos", "Paris", "Rio de Janeiro",
         "San Juan", "Toronto")
city_temps <- data.frame(name = city, temperature = temp)
```

Use the `rank` function to determine the population rank of each state from smallest population size to biggest. Save these ranks in an object called `ranks`, then create a data frame with the state name and its rank. Call the data frame `my_df`.

6. Repeat the previous exercise, but this time order `my_df` so that the states are ordered from least populous to most populous. Hint: create an object `ind` that stores the indexes needed to order the population values. Then use the bracket operator `[` to re-order each column in the data frame.

7. The `na_example` vector represents a series of counts. You can quickly examine the object using:

```
data("na_example")
str(na_example)
#> int [1:1000] 2 1 3 2 1 3 1 4 3 2 ...
```

However, when we compute the average with the function `mean`, we obtain an NA:

```
mean(na_example)
#> [1] NA
```

The `is.na` function returns a logical vector that tells us which entries are NA. Assign this logical vector to an object called `ind` and determine how many NAs does `na_example` have.

8. Now compute the average again, but only for the entries that are not NA. Hint: remember the `!` operator.

---

## 2.11 Vector arithmetics

California had the most murders, but does this mean it is the most dangerous state? What if it just has many more people than any other state? We can quickly confirm that California indeed has the largest population:

```
library(dslabs)
data("murders")
murders$state[which.max(murders$population)]
#> [1] "California"
```

with over 37 million inhabitants. It is therefore unfair to compare the totals if we are interested in learning how safe the state is. What we really should be computing is the

murders per capita. The reports we describe in the motivating section used murders per 100,000 as the unit. To compute this quantity, the powerful vector arithmetic capabilities of R come in handy.

### 2.11.1 Rescaling a vector

In R, arithmetic operations on vectors occur *element-wise*. For a quick example, suppose we have height in inches:

```
inches <- c(69, 62, 66, 70, 70, 73, 67, 73, 67, 70)
```

and want to convert to centimeters. Notice what happens when we multiply `inches` by 2.54:

```
inches * 2.54
#> [1] 175 157 168 178 178 185 170 185 170 178
```

In the line above, we multiplied each element by 2.54. Similarly, if for each entry we want to compute how many inches taller or shorter than 69 inches, the average height for males, we can subtract it from every entry like this:

```
inches - 69
#> [1] 0 -7 -3 1 1 4 -2 4 -2 1
```

### 2.11.2 Two vectors

If we have two vectors of the same length, and we sum them in R, they will be added entry by entry as follows:

$$\begin{pmatrix} a \\ b \\ c \\ d \end{pmatrix} + \begin{pmatrix} e \\ f \\ g \\ h \end{pmatrix} = \begin{pmatrix} a + e \\ b + f \\ c + g \\ d + h \end{pmatrix}$$

The same holds for other mathematical operations, such as `-`, `*` and `/`.

This implies that to compute the murder rates we can simply type:

```
murder_rate <- murders$total / murders$population * 100000
```

Once we do this, we notice that California is no longer near the top of the list. In fact, we can use what we have learned to order the states by murder rate:

```
murders$abb[order(murder_rate)]
#> [1] "VT" "NH" "HI" "ND" "IA" "ID" "UT" "ME" "WY" "OR" "SD" "MN" "MT"
#> [14] "CO" "WA" "WV" "RI" "WI" "NE" "MA" "IN" "KS" "NY" "KY" "AK" "OH"
#> [27] "CT" "NJ" "AL" "IL" "OK" "NC" "NV" "VA" "AR" "TX" "NM" "CA" "FL"
#> [40] "TN" "PA" "AZ" "GA" "MS" "MI" "DE" "SC" "MD" "MO" "LA" "DC"
```

## 2.12 Exercises

1. Previously we created this data frame:

```
temp <- c(35, 88, 42, 84, 81, 30)
city <- c("Beijing", "Lagos", "Paris", "Rio de Janeiro",
         "San Juan", "Toronto")
city_temps <- data.frame(name = city, temperature = temp)
```

Remake the data frame using the code above, but add a line that converts the temperature from Fahrenheit to Celsius. The conversion is  $C = \frac{5}{9} \times (F - 32)$ .

2. What is the following sum  $1 + 1/2^2 + 1/3^2 + \dots 1/100^2$ ? Hint: thanks to Euler, we know it should be close to  $\pi^2/6$ .

3. Compute the per 100,000 murder rate for each state and store it in the object `murder_rate`. Then compute the average murder rate for the US using the function `mean`. What is the average?

## 2.13 Indexing

R provides a powerful and convenient way of indexing vectors. We can, for example, subset a vector based on properties of another vector. In this section, we continue working with our US murders example, which we can load like this:

```
library(dslabs)
data("murders")
```

### 2.13.1 Subsetting with logicals

We have now calculated the murder rate using:

```
murder_rate <- murders$total / murders$population * 100000
```

Imagine you are moving from Italy where, according to an ABC news report, the murder rate is only 0.71 per 100,000. You would prefer to move to a state with a similar murder rate. Another powerful feature of R is that we can use logicals to index vectors. If we compare a vector to a single number, it actually performs the test for each entry. The following is an example related to the question above:

```
ind <- murder_rate < 0.71
```

If we instead want to know if a value is less or equal, we can use:

```
ind <- murder_rate <= 0.71
```

Note that we get back a logical vector with `TRUE` for each entry smaller than or equal to 0.71. To see which states these are, we can leverage the fact that vectors can be indexed with logicals.

```
murders$state[ind]
#> [1] "Hawaii"      "Iowa"          "New Hampshire" "North Dakota"
#> [5] "Vermont"
```

In order to count how many are `TRUE`, the function `sum` returns the sum of the entries of a vector and logical vectors get *coerced* to numeric with `TRUE` coded as 1 and `FALSE` as 0. Thus we can count the states using:

```
sum(ind)
#> [1] 5
```

### 2.13.2 Logical operators

Suppose we like the mountains and we want to move to a safe state in the western region of the country. We want the murder rate to be at most 1. In this case, we want two different things to be true. Here we can use the logical operator *and*, which in R is represented with `&`. This operation results in `TRUE` only when both logicals are `TRUE`. To see this, consider this example:

```
TRUE & TRUE
#> [1] TRUE
TRUE & FALSE
#> [1] FALSE
FALSE & FALSE
#> [1] FALSE
```

For our example, we can form two logicals:

```
west <- murders$region == "West"
safe <- murder_rate <= 1
```

and we can use the `&` to get a vector of logicals that tells us which states satisfy both conditions:

```
ind <- safe & west
murders$state[ind]
#> [1] "Hawaii" "Idaho"  "Oregon" "Utah"   "Wyoming"
```

### 2.13.3 which

Suppose we want to look up California's murder rate. For this type of operation, it is convenient to convert vectors of logicals into indexes instead of keeping long vectors of

logicals. The function `which` tells us which entries of a logical vector are TRUE. So we can type:

```
ind <- which(murders$state == "California")
murder_rate[ind]
#> [1] 3.37
```

### 2.13.4 `match`

If instead of just one state we want to find out the murder rates for several states, say New York, Florida, and Texas, we can use the function `match`. This function tells us which indexes of a second vector match each of the entries of a first vector:

```
ind <- match(c("New York", "Florida", "Texas"), murders$state)
ind
#> [1] 33 10 44
```

Now we can look at the murder rates:

```
murder_rate[ind]
#> [1] 2.67 3.40 3.20
```

### 2.13.5 `%in%`

If rather than an index we want a logical that tells us whether or not each element of a first vector is in a second, we can use the function `%in%`. Let's imagine you are not sure if Boston, Dakota, and Washington are states. You can find out like this:

```
c("Boston", "Dakota", "Washington") %in% murders$state
#> [1] FALSE FALSE TRUE
```

Note that we will be using `%in%` often throughout the book.

**Advanced:** There is a connection between `match` and `%in%` through `which`. To see this, notice that the following two lines produce the same index (although in different order):

```
match(c("New York", "Florida", "Texas"), murders$state)
#> [1] 33 10 44
which(murders$state %in% c("New York", "Florida", "Texas"))
#> [1] 10 33 44
```

---

## 2.14 Exercises

Start by loading the library and data.

```
library(dslabs)
data(murders)
```

1. Compute the per 100,000 murder rate for each state and store it in an object called `murder_rate`. Then use logical operators to create a logical vector named `low` that tells us which entries of `murder_rate` are lower than 1.
2. Now use the results from the previous exercise and the function `which` to determine the indices of `murder_rate` associated with values lower than 1.
3. Use the results from the previous exercise to report the names of the states with murder rates lower than 1.
4. Now extend the code from exercises 2 and 3 to report the states in the Northeast with murder rates lower than 1. Hint: use the previously defined logical vector `low` and the logical operator `&`.
5. In a previous exercise we computed the murder rate for each state and the average of these numbers. How many states are below the average?
6. Use the `match` function to identify the states with abbreviations AK, MI, and IA. Hint: start by defining an index of the entries of `murders$abb` that match the three abbreviations, then use the `[]` operator to extract the states.
7. Use the `%in%` operator to create a logical vector that answers the question: which of the following are actual abbreviations: MA, ME, MI, MO, MU?
8. Extend the code you used in exercise 7 to report the one entry that is **not** an actual abbreviation. Hint: use the `!` operator, which turns `FALSE` into `TRUE` and vice versa, then `which` to obtain an index.

---

## 2.15 Basic plots

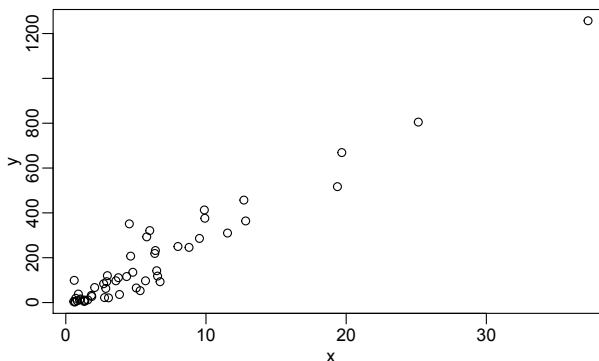
In Chapter 7 we describe an add-on package that provides a powerful approach to producing plots in R. We then have an entire part on Data Visualization in which we provide many examples. Here we briefly describe some of the functions that are available in a basic R installation.

### 2.15.1 plot

The `plot` function can be used to make scatterplots. Here is a plot of total murders versus population.

```
x <- murders$population / 10^6
y <- murders$total
plot(x, y)
```





For a quick plot that avoids accessing variables twice, we can use the `with` function:

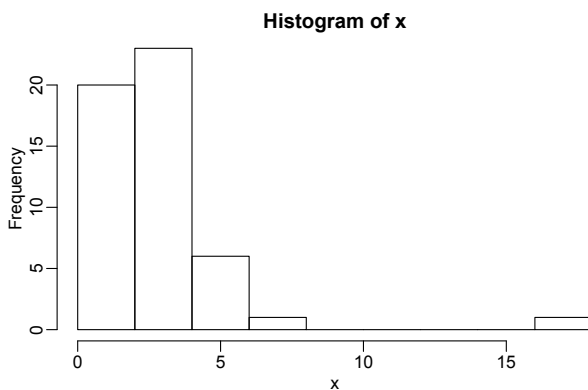
```
with(murders, plot(population, total))
```

The function `with` lets us use the `murders` column names in the `plot` function. It also works with any data frames and any function.

### 2.15.2 hist

We will describe histograms as they relate to distributions in the Data Visualization part of the book. Here we will simply note that histograms are a powerful graphical summary of a list of numbers that gives you a general overview of the types of values you have. We can make a histogram of our murder rates by simply typing:

```
x <- with(murders, total / population * 100000)
hist(x)
```



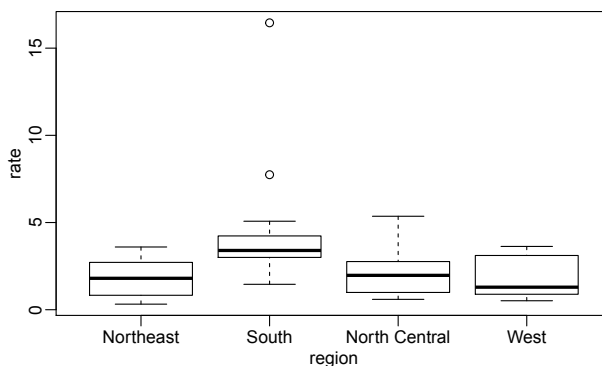
We can see that there is a wide range of values with most of them between 2 and 3 and one very extreme case with a murder rate of more than 15:

```
murders$state[which.max(x)]
#> [1] "District of Columbia"
```

### 2.15.3 boxplot

Boxplots will also be described in the Data Visualization part of the book. They provide a more terse summary than histograms, but they are easier to stack with other boxplots. For example, here we can use them to compare the different regions:

```
murders$rate <- with(murders, total / population * 100000)
boxplot(rate~region, data = murders)
```

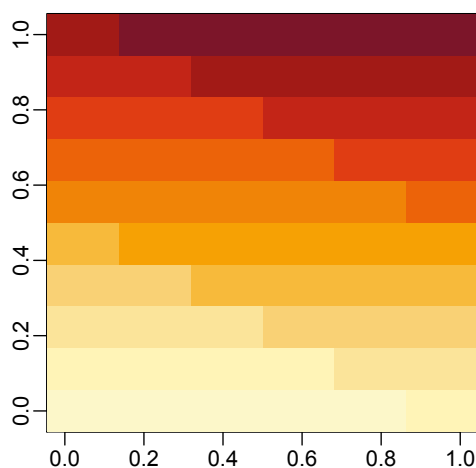


We can see that the South has higher murder rates than the other three regions.

### 2.15.4 image

The image function displays the values in a matrix using color. Here is a quick example:

```
x <- matrix(1:120, 12, 10)
image(x)
```



---

## 2.16 Exercises

1. We made a plot of total murders versus population and noted a strong relationship. Not surprisingly, states with larger populations had more murders.

```
library(dslabs)
data(murders)
population_in_millions <- murders$population/10^6
total_gun_murders <- murders$total
plot(population_in_millions, total_gun_murders)
```

Keep in mind that many states have populations below 5 million and are bunched up. We may gain further insights from making this plot in the log scale. Transform the variables using the `log10` transformation and then plot them.

2. Create a histogram of the state populations.
3. Generate boxplots of the state populations by region.