



# Data Insight Foundations

Step-by-Step Data Analysis with R

—

Nikita Tkachenko

Apress®

# Data Insight Foundations

Nikita Tkachenko

# Data Insight Foundations

Step-by-Step Data Analysis with R

Apress®

Nikita Tkachenko  
San Francisco, CA, USA

ISBN-13 (pbk): 979-8-8688-0579-0

ISBN-13 (electronic): 979-8-8688-0580-6

<https://doi.org/10.1007/979-8-8688-0580-6>

Copyright © 2025 by Nikita Tkachenko

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

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

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

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

Managing Director, Apress Media LLC: Welmoed Spahr

Acquisitions Editor: Shaul Elson

Development Editor: Laura Berendson

Editorial Assistant: Gryffin Winkler

Cover designed by eStudioCalamar

Cover image designed by Freepik ([www.freepik.com](http://www.freepik.com))

Distributed to the book trade worldwide by Springer Science+Business Media LLC, 1 New York Plaza, Suite 4600, New York, NY 10004. Phone 1-800-SPRINGER, fax (201) 348-4505, e-mail [orders-ny@springer-sbm.com](mailto:orders-ny@springer-sbm.com), or visit [www.springeronline.com](http://www.springeronline.com). Apress Media, LLC is a California LLC and the sole member (owner) is Springer Science + Business Media Finance Inc (SSBM Finance Inc). SSBM Finance Inc is a **Delaware** corporation.

For information on translations, please e-mail [booktranslations@springernature.com](mailto:booktranslations@springernature.com); for reprint, paperback, or audio rights, please e-mail [bookpermissions@springernature.com](mailto:bookpermissions@springernature.com).

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

Any source code or other supplementary material referenced by the author in this book is available to readers on the Github repository. For more detailed information, please visit <https://www.apress.com/gp/services/source-code>.

If disposing of this product, please recycle the paper

*I dedicate this book to my late grandfather, Vladimir  
Zamashikov.*



# Contents

About the Author .....	xiii
About the Technical Reviewer .....	xv
Acknowledgements .....	xvii
Introduction .....	xix
Setting Up R and RStudio .....	xxi
 <b>Part I Working with Data</b>	
<b>1 Data Manipulation .....</b>	<b>3</b>
1.1 Basics .....	3
1.1.1 Data Types .....	5
1.2 Downloading Data .....	7
1.2.1 Example Data .....	8
1.3 Basic Data Management with dplyr .....	9
1.3.1 select() .....	9
1.3.2 filter() .....	10
1.3.3 arrange() .....	11
1.3.4 mutate() .....	12
1.3.5 case_match() .....	13
1.3.6 summarize() .....	13
1.3.7 group_by() .....	14
1.3.8 ungroup() .....	15
1.3.9 .by .....	15
1.3.10 rowwise() .....	16
1.3.11 count() .....	17
1.3.12 rename() .....	18
1.3.13 row_number() .....	18
1.3.14 skim() .....	18
1.4 Exploring Date and Time with lubridate .....	19
1.4.1 ymd(), md(), hms(), ymd_hms() .....	19
1.4.2 year(), month(), day() .....	20
1.5 Summary .....	21

<b>2</b>	<b>Tidy Data</b>	23
2.1	Example	24
2.2	<code>tidyr</code>	25
2.2.1	<code>pivot_longer()</code>	25
2.2.2	<code>pivot_wider()</code>	26
2.2.3	<code>separate()</code> and <code>unite()</code>	27
2.3	<code>tibble()</code> and <code>tribble()</code>	29
2.4	janitor: Clean Your Data	29
2.4.1	<code>clean_names()</code>	30
2.4.2	<code>remove_empty()</code>	30
2.4.3	<code>remove_constant()</code>	31
2.4.4	<code>convert_to_date()</code> and <code>convert_to_datetime()</code>	32
2.4.5	<code>row_to_names()</code>	32
2.5	Summary	33
<b>3</b>	<b>Relational Data</b>	35
3.1	Relationship Types	35
3.1.1	One to One (1:1)	35
3.1.2	One to Many (1:M)	36
3.1.3	Many to Many (M:N)	36
3.2	The Concept of Keys	36
3.3	Types of Joins	37
3.3.1	Outer Joins	37
3.3.2	Filtering Joins	40
3.4	Visualizing Relationships	41
3.5	Summary	42
<b>4</b>	<b>Data Validation</b>	43
4.1	Manual Inspection	43
4.2	Handling Data Issues	43
4.2.1	Assert Your Conditions	44
4.2.2	Precise Validation with Pointblank	45
4.3	Summary	47
<b>5</b>	<b>Imputation</b>	49
5.1	Types of Missing Data	49
5.2	Dealing with Missing Data	49
5.2.1	Explicitly Handling Missing Data with <code>complete()</code>	50
5.2.2	Simple Imputations	52
5.2.3	Best Worst-Case and Worst Best-Case Scenarios	56
5.2.4	Multiple Imputations	56
5.3	Summary	59
5.3.1	Table of Imputations	60
<b>Part II Reproducible Research</b>		
<b>6</b>	<b>Reproducible Research</b>	65
6.1	Literate Programming	66
6.2	Summary	67
<b>7</b>	<b>Reproducible Environment</b>	69
7.1	<code>renv</code>	69



7.1.1	Workflow .....	70
7.2	Computational Environments .....	70
7.3	Summary .....	71
<b>8</b>	<b>Introduction to Command Line .....</b>	<b>73</b>
8.1	Learning Basic Commands .....	73
8.2	Getting Started with Nano .....	74
<b>9</b>	<b>Version Control with Git and Github .....</b>	<b>75</b>
9.1	Git and GitHub .....	76
9.2	Basics .....	76
9.3	Guide to Using .gitignore .....	78
9.3.1	Specifying Files to Ignore .....	78
9.3.2	.gitignore in Other Programs .....	79
9.4	Summary .....	79
<b>10</b>	<b>Style and Lint Your Code .....</b>	<b>81</b>
10.1	Tidyverse Style Guide .....	81
10.1.1	White Spaces and Indentation .....	81
10.1.2	Naming Conventions .....	82
10.1.3	Braces .....	82
10.1.4	Comments .....	83
10.1.5	Long Functions .....	83
10.2	Formatter .....	84
10.3	Linters .....	84
10.4	Summary .....	85
<b>11</b>	<b>Modular Code .....</b>	<b>87</b>
11.1	Reuse Functions .....	87
11.2	Split It .....	88
11.3	box It .....	89
11.4	Summary .....	91
 <b>Part III Literature Review and Writing</b>		
<b>12</b>	<b>Literature Review .....</b>	<b>95</b>
12.1	Search .....	95
12.2	Reference Management .....	96
12.3	Reading .....	96
12.4	Note Taking .....	97
12.5	Summary .....	97
<b>13</b>	<b>Write .....</b>	<b>99</b>
13.1	WYSIWYG .....	99
13.2	Markup Languages .....	99
13.2.1	HTML .....	100
13.2.2	LaTeX .....	100
13.2.3	Markdown .....	100
13.2.4	Yet Another Markup Language (YAML) .....	101
13.2.5	Pandoc .....	101
13.3	Quarto .....	101
13.3.1	Your First Document .....	102
13.4	Summary .....	104

<b>14</b>	<b>Layout and References</b>	105
14.1	Knitr	105
14.2	Div Blocks	106
14.3	Diagrams	108
14.4	Citations	108
14.5	Summary	109
<b>15</b>	<b>Collaboration and Templating</b>	111
15.1	Streamlining Collaboration with <code>trackdown</code>	111
15.2	Templating	112
15.3	Summary	113
 <b>Part IV Collecting the Data</b>		
<b>16</b>	<b>Total Survey Error (TSE)</b>	117
16.1	Representation—The People You Ask	119
16.1.1	Sampling	120
16.1.2	What Responses Depend On	122
16.2	Question Design	123
16.2.1	Answering Questions	123
16.2.2	Guidelines for Effective Question Design	124
16.2.3	The Impact of Question Order	125
16.2.4	Survey Says: Bacon!	125
16.2.5	Types of Question Formats	126
16.2.6	Likert Scales	127
16.2.7	“I don’t know” and N/A Options	127
16.2.8	Survey Length	128
16.2.9	Survey Invitation	128
16.2.10	Iterative Design	129
16.3	Survey Tools	129
16.3.1	Physical Survey Methods	129
16.3.2	Digital Survey Tools	130
16.3.3	Participant Recruitment Platforms	131
16.4	Summary	131
<b>17</b>	<b>Document</b>	133
17.1	Principles of Documentation	133
17.1.1	Hierarchical Structure of Documentation	134
17.2	Physical and Electronic Documentation	135
17.3	Data Organization in Spreadsheets	136
17.3.1	Spreadsheet Organization	136
17.3.2	Data Input	138
17.4	Administration	139
17.5	Accounting in Research	139
17.6	Communication	140
17.7	Summary	141
<b>18</b>	<b>APIs</b>	143
18.1	API Basics	143
18.2	Utilizing APIs in R	145
18.3	Qualtrics API	147
18.4	Integrating Google Services with R	148

18.5	OpenAI's API .....	150
18.6	Summary .....	152

## Part V Presenting the Data

<b>19</b>	<b>Data Visualization Fundamentals</b> .....	155
19.1	Perceptual Processing .....	155
19.1.1	Preattentive Processing .....	156
19.2	Visual Encoding .....	161
19.2.1	Evaluating Graphs .....	162
19.3	Summary .....	163
<b>20</b>	<b>Data Visualization</b> .....	165
20.1	Exploratory vs. Explanatory .....	165
20.2	The Grammar of Graphics .....	166
20.2.1	Graphical Interfaces for <code>ggplot2</code> .....	173
20.3	Interactive Plots .....	174
20.3.1	HTML Widgets .....	175
20.4	Summary .....	177
<b>21</b>	<b>A Graph for the Job</b> .....	179
21.1	Category Comparison .....	179
21.1.1	Lollipop Chart .....	180
21.1.2	Bullet Graph .....	180
21.2	Distribution .....	181
21.2.1	Density Plot .....	182
21.2.2	Frequency Polygon .....	182
21.2.3	Box Plot .....	183
21.2.4	Violin Plot .....	183
21.2.5	Bee Hive Plot .....	184
21.2.6	Rain Cloud Plot .....	185
21.2.7	Margins .....	186
21.3	Proportions .....	186
21.3.1	Stacked Bar Charts .....	186
21.3.2	Pie Chart .....	188
21.3.3	Waffle Chart .....	189
21.3.4	Treemaps .....	189
21.4	Correlation .....	189
21.4.1	Scatter Plot .....	190
21.4.2	Correlograms .....	190
21.5	Change Over Time .....	191
21.5.1	Line Chart .....	191
21.5.2	Waterfall Chart .....	192
21.6	Summary .....	193
<b>22</b>	<b>Color Data</b> .....	195
22.1	What Colors to Choose .....	195
22.1.1	Complementary Harmony with a Positive/Negative Connotation .....	195
22.1.2	Near Complementary Harmony for Highlighting Two Series Where One Is the Primary Focus .....	196
22.1.3	Analogous/Triadic Harmony for Highlighting Three Series .....	197
22.1.4	Highlighting One Series Against Two Related Series .....	198

22.1.5	Analogous Complementary for One Main Series and Its Three Secondary .....	199
22.1.6	Double Complementary for Two Pairs Where One Pair Is Dominant .....	199
22.1.7	Rectangular or Square Complementary Scheme for Four Series of Equal Emphasis .....	200
22.1.8	Sequential .....	200
22.1.9	Divergent .....	201
22.1.10	Prebuilt .....	202
22.2	Color Systems .....	203
22.2.1	Warning: Colormaps Might Increase Risk of Death! .....	205
22.2.2	So What Should You Use? .....	205
22.3	Summary .....	207
<b>23</b>	<b>Make Tables .....</b>	<b>209</b>
23.1	gt Tables .....	209
23.1.1	Prepare Your Data .....	210
23.2	DT Tables .....	217
23.3	Summary .....	218
<b>Epilogue .....</b>		<b>219</b>
<b>References .....</b>		<b>221</b>
<b>Index .....</b>		<b>223</b>

## About the Author



**Nikita Tkachenko** leads Evalyn, a consulting agency specializing in AI-driven customer service audits and data analytics solutions. He helps organizations of all sizes harness AI and data to optimize decision-making, streamline operations, and enhance customer experiences. With a strong foundation in research and analytics, Nikita also teaches courses on research tools, mentors students, and conducts academic research at the University of San Francisco.

## About the Technical Reviewer



**Sijo Valayakkad Manikandan** is an accomplished artificial intelligence and data science leader. He has extensive experience managing large-scale data science projects for Fortune 500 corporations. Sijo holds a Master of Science degree in Business Analytics from the renowned University of Texas at Austin and a Bachelor of Engineering degree from the Birla Institute of Technology and Science in Pilani, India. He combines his academic excellence and practical expertise to deliver exceptional results.

Sijo has made significant contributions to the field of data science and research, not only through his professional and academic achievements but also through his dedication to the community.

He is a member of several distinguished organizations, such as the American Statistical Association, and has conducted independent research. Sijo has also actively reviewed academic and professional books, mentored junior data scientists, and guided early-stage startups. Moreover, his expertise has led him to serve on the jury of several prestigious awards, including The Webby Awards.

As a dedicated researcher and thought leader, Sijo continues to profoundly impact the field of data science, inspiring a new generation of data scientists. His contributions have advanced the field of data science and research and helped shape the future of this rapidly evolving industry, making him an invaluable asset to the community and a visionary in his field.

# Acknowledgements

I dedicate this book to my late grandfather, Vladimir Zamashikov.

I extend my deepest gratitude to Thomas Weinandy, who convinced me to transform my rough draft into a published book and guided me through the process. Without him, this book would not have been possible.

I am deeply indebted to Parsa Rahimi for his unwavering support, meticulous review of all my chapters, edits, and ideas, and for providing artworks featuring SF.

My thanks also go to Alessandra Cassar and Michael Jonas for their support throughout my academic journey. Special thanks go to Peter Lorentzen, Ker Gibbs, and Jesse Anttila-Hughes for their invaluable guidance and support.

I appreciate Shivani Shukla, Bruce Wydick, Misha Gipsman, Mario Lim, Steve Trettel, Mehmet Emre, Andrew Hobbs, Arman Khachiyani, Robizon Khubulashvili, and Konrad Posch for reviewing chapters and providing their support. Additionally, I am grateful to the faculty, librarians, writing center staff, and the Center for Business Studies and Innovation (CBSI) at the University of San Francisco for their guidance.

Gordon Getty deserves thanks for inspiring me to write more confidently about my skills.

The team at APress, especially Shaul Elson, Commissioning Editor, Laura Berendson, Development Editor, and Gryffin Winkler, Editorial Assistant, guided me through the publishing process. My technical reviewer, Sijo Manikandan, provided comments that significantly enriched the book with additional examples and explanations. Special thanks go to everyone at APress whose efforts made this book a reality.

I gratefully acknowledge John Chetwynd, Jake Consgrrove, Timofei Lopukhov, and Alexander Rom, whose support, curiosity, and confidence in my work are deeply appreciated.

I must acknowledge my family, including my mother Ekaterina Tkachenko, whose belief in me inspired me to launch this project, my father Anton Tkachenko, stepmother Daria Tkachenko, and grandparents Valentina and Alexander Trachenko and Raisa Zamashikova.

A special note of thanks goes to Anastasia Ternovskaya for her unwavering belief in me.

I very much appreciate Maria Aksutenko for making the artwork of the “researcher-cat.”

Special thanks also go to the Posit team for their invaluable contributions and support of the community.

Finally, my heartfelt appreciation goes to the open-source community and everyone who has generously shared their knowledge online. You made this possible, and I am thrilled to contribute to the conversation. I hope you find this book useful in your journey toward mastering research.

# Introduction

This book was born from my frustrations and experiences in higher education and professional work. It originated from notes and materials from a Spring 2023 course in survey design, inspired by the enthusiastic response and insightful questions from my students.

Young data professionals typically learn about models, experiments, and theories during their classes, frequently returning to that knowledge. However, executing high-quality research and analysis requires a deeper understanding of the tools and the "how" rather than just the "what" and "why." This knowledge often transcends what is taught within the constraints of a standard curriculum. In this book, I aim to bridge that gap, helping you move from knowing what you want to do to understanding how to do it. I have distilled hundreds of hours of frustration into these chapters, so you won't have to traverse that path yourself.

This book is not a comprehensive guide; if that's what you're seeking, you may want to look elsewhere. Instead, the book can serve as a map, outlining the necessary tools and topics for your research journey. The goal is to build your intuition and provide pointers for where to find more detailed information. The chapters are deliberately concise and to the point, aiming to reveal and enlighten rather than bore. You'll learn about efficient data management, reproducible research, literature review and writing practices, and effective data visualization.

Initially inspired by my journey through graduate school in economics, this book offers value across disciplines. It contains essential insights for anyone engaged in data-related work, from the humanities to data analytics and the sciences. Whether you are refining your expertise or new to data analytics, this book promises to offer something of value.

Examples provided are primarily in R, making a basic understanding of the language advantageous but not essential. Several chapters, especially those focusing on theory, require no programming knowledge at all. A diverse audience, including web developers, mathematicians, data analysts, and economists, has found the material beneficial. The book is designed to be inclusive, offering insights irrespective of your programming proficiency or professional background.

Its structure allows for flexible reading paths; you may explore the chapters in sequence for a systematic learning experience or navigate directly to the topics most relevant to you.



# Setting Up R and RStudio

Welcome to the exciting world of data analysis with R, a language crafted specifically for statistical analysis and data visualization. R's user-friendly syntax and reproducibility make it an ideal choice for both novices and professionals. However, before diving into data exploration and modeling, it's essential to differentiate between R, the programming language, and RStudio, the integrated development environment (IDE) that enhances R's functionality.

## Download and Install R

R is maintained and distributed through the Comprehensive R Archive Network (CRAN), ensuring your access to the latest version and resources.

### *For macOS Users:*

1. Navigate to the [CRAN website](#).
2. Click on "Download R for macOS."
3. Select the appropriate version:
  - For Apple Silicon (e.g., M1, M2), download the version with "-arm64" (e.g., R-4.2.2-arm64.pkg) in its name.
  - For Intel-based Macs, select the version without "-arm64" (e.g., R-4.2.2.pkg).
4. Follow the installation wizard. The default settings are typically sufficient.

### *For Windows Users:*

1. Visit the [CRAN website](#).
2. Choose "Download R for Windows."
3. Select "base" and then the first link at the top of the page (e.g., Download R-4.2.2 for Windows).
4. The installer will guide you through the process. Stick with the default settings for a smooth installation.

5. Additionally, Windows users should download Rtools, which are crucial for compiling packages from sources. Visit [Rtools](#), match the Rtools version with your R version, and follow the installer instructions.

## Download and Install RStudio

RStudio provides a user-friendly interface for R, akin to what Microsoft Word offers for text but tailored to R scripting.

- To download RStudio, head to the [RStudio download page](#).
- Click on “DOWNLOAD RSTUDIO DESKTOP FOR WINDOWS” or select your operating system for detailed instructions.

## Configure RStudio

Enhance your RStudio experience with these initial setup tips:

- **Change Theme:** Shift from the default theme to a dark theme for improved readability. Go to Tools > Global Options > Appearance, and select “Dracula.” Click “Apply.”
- **Install Fira Code Font:** For a modern coding aesthetic, install the “Fira Code” font, which supports programming ligatures. Instructions can be found at [Fira Code on GitHub](#). After installation, apply this font in RStudio under Appearance.

## Install Packages

Packages extend R’s functionality. Install them easily with commands in the RStudio console:

```
# Install a single package
install.packages("tidyverse")

# Install multiple packages
install.packages(c("tidyverse", "gapminder"))
```

To use installed packages, load them into your session:

```
library(tidyverse)
```

Now that we have R and RStudio up and running, let’s dive into some fundamental data manipulation techniques in R.

# **Part I**

## **Working with Data**

# Chapter 1

## Data Manipulation

In data analysis, visualization and manipulation are essential for understanding and communicating complex information. R, a powerful programming language for data analysis, offers a variety of packages that enable the creation of visually compelling plots and the streamlining of data manipulation. One of the most user-friendly and widely used collections of R packages<sup>1</sup> is tidyverse, developed by Hadley Wickham, chief scientist at Posit (RStudio). Tidyverse includes packages that cover all common tasks and can be installed with `install.packages("tidyverse")` and activated using `library("tidyverse")`. In this introduction, we will cover the basics of tidyverse using `readr` for data reading, `dplyr` for data manipulation, `tidyr` for data tidying, and, later in the book, `ggplot2` for data visualization. For more information about tidyverse, visit their website <https://www.tidyverse.org/>.

### 1.1 Basics

Let's kick off with some fundamental concepts! R can be employed as a simple calculator.

```
# A "#" is used to annotate comments!  
2 + 2
```

```
#> [1] 4
```

```
2 * 4
```

```
#> [1] 8
```

---

<sup>1</sup> A package is a collection of prewritten functions, data, and documentation that enhances the capabilities of the R programming language for specific tasks.

```
2^8
```

```
#> [1] 256
```

```
(1 + 3) / (3 + 5)
```

```
#> [1] 0.5
```

```
log(10) # Calculates the natural log of 10!
```

```
#> [1] 2.302585
```

R allows for defining variables and performing operations on them. Both `=` and `<-` can be used for assigning values to a variable name, though `<-` is preferred to avoid confusion and certain errors.

```
x <- 2 # Equivalent to x = 2  
x * 4
```

```
#> [1] 8
```

The command `x <- 2` assigns the value 2 to `x`. Thus, when we subsequently type `x * 4`, R replaces `x` with 2 to evaluate `2 * 4` and obtain 8. The value of `x` can be updated as needed using `=` or `<-`. Bear in mind that R is case sensitive, so `X` and `x` are recognized as different variables.

```
x
```

```
#> [1] 2
```

```
(x <- x * 5) # Wrapping with (...) prints the variable
```

```
#> [1] 10
```

To further explore operations in R, the following table presents a comprehensive overview of basic arithmetic, comparison, and logical operators you might need.

Operator	Description	Example	Result
+	Addition	3 + 2	5
-	Subtraction	5 - 2	3
*	Multiplication	3 * 2	6
/	Division	6 / 2	3
^	Exponentiation	2 ^ 3	8
%%	Modulus (remainder)	5 %% 2	1
%/%	Integer division	5 %/% 2	2
==	Equal to	2 == 2	TRUE
!=	Not equal to	2 != 3	TRUE
<	Less than	2 < 3	TRUE
>	Greater than	3 > 2	TRUE
<=	Less than or equal to	2 <= 2	TRUE
>=	Greater than or equal to	2 >= 2	TRUE
&	Logical AND	TRUE & FALSE	FALSE
	Logical OR	TRUE   FALSE	TRUE
!	Logical NOT	!TRUE	FALSE

### 1.1.1 Data Types

R possesses a multitude of data types and classes, including `data.frames`, which are akin to Excel spreadsheets with columns and rows. Initially, we'll examine vectors. Vectors can store multiple values of the same type, with the most basic ones being numeric, character, and logical.

```
x
```

```
#> [1] 10
```

```
class(x)
```

```
#> [1] "numeric"
```

```
(true_or_false <- TRUE)
```

```
#> [1] TRUE
```

```
class(true_or_false)
```

```
#> [1] "logical"
```

```
(name <- "Parsa Rahimi")
```