

Introduction to R

Data analysis and statistical computing for quantitative social science research

Laboratorio de Investigación para el Desarrollo del Ecuador

Instructor: Daniel Sánchez, MA

Module length: 15 hours

Course level: Intermediate

Prerequisite knowledge: Basic computer skills, Git, shell.

Corequisite knowledge: **Introduction to Statistics** module.

GitHub repository: <https://github.com/laboratoriolide/intro-to-r>

1 Module overview

This module introduces the R programming language basics for data analysis and statistical computing. This is a companion module to the *Introduction to Statistics* module, which provides the foundations of statistics for social science research. The material is designed for those who aim to use R as a tool for quantitative research methods, such as statistics, econometrics, and (to a lesser degree), data science/machine learning.

I assume little prior coding experience: familiarity with basic computer skills, Git and the command line. Further, the course will be relatively fast-paced, combining self-study with the DataCamp platform and live coding. We will cover many important concepts in the language, however, the focus is not on technical programming but rather on understanding the language as a tool for statistical analysis.

2 Module contents

The following is a planned outline of the course. This may change depending on the pace of the class. It is sometimes expected that students complete tasks before each lecture to make the most of the class time.

- **Lecture 1:** Intro
 - Introduction to the R language: brief history, advantages and applications
 - The development environment: R, RStudio, IDEs, and packages.
 - The RStudio interface: R scripts, file viewer, console, and environment.
 - Ensuring reproducibility: filepaths, working adirectories, and project organization; using R projects with RStudio.
 - Installing, loading and updating packages, best practices.
 - Basic R syntax: variables, assigners, data types, structures, basic and logical operators.
 - Getting help with R: `?`, `help()`. Using R documentation, CRAN, StackOverflow & other resources.
 - Data frames: definition and basic manipulation with base R syntax: `head()`, `tail()`, `str()`, `summary()`. Basic R graphics.
 - Artificial intelligence tools for R coding: ChatGPT, GitHub Copilot, Microsoft Copilot.
 - Updating R and R packages.
 - **Complete before class:** installation of R and RStudio, cloning of the course's GitHub repository and completion of the first DataCamp course, *Introduction to R for Finance*.
- **Lecture 2:** Introduction to the tidyverse and importing data
 - Base R vs the *tidyverse*
 - Introduction to the *tidyverse*: the tidy process, tidy data, the pipe (`%>%`) operator
 - Importing data with base R functions
 - Importing data with `readr`, `haven` from different file formats
 - Importing Excel files and preliminary manipulation with `readxl`
 - Downloading data from R packages, GitHub and the web/URLs
 - Tibbles: the tidy data frame; using `glimpse()` to understand data
 - dplyr: Selecting and renaming columns, pipe workflows, `transmute`.
 - Complete before class: DataCamp course *Introduction to importing data in R* and *Introduction to the Tidyverse* (select chapters) and make sure you have access to the module's `data` folder in the GitHub repository.
- **Lecture 3:** More on data manipulation
 - Filtering rows, sorting, grouping, and summarizing data.

- Using `mutate()` to create new columns: `if_else`, `case_when`, `case_match()`.
- Joins, semijoins and anti joins: applications for researchers.
- The `bind_rows()` function for appending data frames.
- The `across()` function for applying functions to multiple columns.
- **Lecture 4:** Data visualization with ggplot2
 - Introduction to the `ggplot2` package: the grammar of graphics.
 - The `ggplot()` function: aesthetics, geoms, and layers.
 - Scatter plots, line plots, bar plots, histograms, and box plots.
 - Faceting: `facet_wrap()` and `facet_grid()`.
 - Themes and customization: `theme()`, `labs()`, `scale_x_continuous()`, `scale_fill_manual()`.
 - Saving plots, changing plot size, and exporting to different formats.
- **Lecture 5:** Advanced topics in data cleaning
 - Understanding long vs. wide data formats; reshaping data with `tidyr`.
 - The `gather()` and `spread()` equivalents to `pivot_longer()` and `pivot_wider()`.
 - The `separate()` and `unite()` functions for splitting and combining columns.
 - The `separate_rows()` function for splitting rows with multiple values.
 - Using the `reshape` package.
 - Expanding data: `expand()` and `complete()` functions.
 - Missing data: identifying, handling, and imputing missing values.
- **Lecture 6:** Reporting and reproducibility
 - The concept of reproducible research.
 - Introduction to R Markdown: the R Markdown document, YAML header, code chunks.
 - Knitting documents: HTML, PDF, Word, and slides.
 - Customizing R Markdown documents: themes, templates, and CSS.
 - Quarto: the next generation of R Markdown.
 - Math formulas: incorporating LaTeX into R Markdown.
 - Time-permitting: `.Rnw` files and Knitr/Sweave.
 - **Complete before class:** DataCamp course *Reporting with R Markdown*. Making sure you have L^AT_EX installed in your computer, or TinyTeX as a package in R¹.
- **Lecture 7:** Select topics in R programming
 - Control structures, conditional statements, and logical operators; relationships to prepackaged functions.
 - Functions: when to use, how to define functions, arguments, return values.
 - The `apply` family of functions: `apply()`, `lapply()`, `sapply()`, `mapply()`.

¹TinyTeX is a lightweight, portable, cross-platform, and easy-to-maintain LaTeX distribution. It is available as a package in R, and can be installed with `tinytex::install_tinytex()`

- Loops: while and for loops, the `map` family of functions.
- **Lecture 8:** Further advanced topics with R
 - Time and date management: `lubridate`.
 - String manipulation with `stringr`.
 - Fuzzy joins: string distance and the `fuzzyjoin` package.
 - Time-permitting: regex, text analysis.
- **Lecture 9:** High-performance computing in R
 - Why high performance? The need for speed in data analysis.
 - Loading very large datasets: `data.table` and `fread()`.
 - Reading many files at once: `purrr` and `map()` with `fread()`.
 - Grouping and summarizing large datasets: `dplyr` and `data.table` comparison.
 - Other alternatives: `bigmemoryr`, `ff`, and SQL databases.
 - Time-permitting: Introduction to parallel computing: the `parallel` package.
- **Lecture 10:** Advanced data access with R
 - Web scraping: the `rvest` package.
 - Accessing twitter data: the `rtweet` package.
 - Accessing financial data: the `quantmod` package.
 - Time-permitting: APIs, applications with open data.
- Other advanced topics:
 - Spatial data analysis with `sf` and `sp`.
 - R Setup with VS Code
 - R in the cloud: RStudio Server, RStudio Cloud, and Google Colab.
 - Shiny apps: building interactive web applications with R.
 - AI chatbot packages: OpenAI tokens, the `chatgpt` and related packages.
 - Using `usethis` for Git version control and R project management.

3 Evaluation

The final grade will be calculated as follows:

Component	Percentage
Attendance	15%
DataCamp assignments	50%
Participation	5%
Task Completion	30%

3.1 Attendance

I do not have any special requirements for attendance other than the requirements set by the program. Consult the program's regulation handbook for more information.

3.2 Assignments

There will be weekly DataCamp assignments that reinforce the material taught in this module as well as that of the companion module *Intro to Statistics* module. These will be either DataCamp courses, projects or exercises, and are graded on a pass/fail basis. Please check the DataCamp platform for the due dates of each assignment.

Apart from the DataCamp assignments, you may need to complete some self-study or autonomous practice work **before** each lecture. This is necessary to make the most of the class time, as programming is a skill that requires practice. Please see above for all tasks to be completed before each lecture.

3.2.1 Tasks

Tasks are small exercises or activities that are due before each lecture. They may be small coding exercises, readings or installations that you **must** complete before class. These will be graded on a pass/fail basis. We may discuss these tasks in class or in the Slack channel, so it is important that you complete them on time.

3.2.2 Tentative DataCamp assignments

This is a list of the DataCamp courses that will be assigned as part of the module. This list may change, so please check the DataCamp platform for the most up-to-date information.

- Introduction to R for Finance
- Introduction to Importing data in R
- Introduction to the Tidyverse
- Data Manipulation with dplyr
- Data Visualization with ggplot2
- Reporting with R Markdown
- Reshaping Data with tidyr (select chapters)
- Intermediate R for Finance

4 Course materials

All course materials will be provided in the course’s GitHub repository. This includes lecture slides, readings, datasets, assignments and any other relevant material. The repository will be updated regularly, so please check it often for new material. I recommend using a Git client, such as GitHub Desktop or GitKraken, to keep your local repository up-to-date.

4.1 Textbook

There is no required textbook for this course, as I will provide slides for all lectures. However, I recommend the following books for those who want to delve deeper into the material. These were used as references for the course.

- *R for Data Science*, Wickham, Çetinkaya-Rundel, and Grolemund (2023): Potentially the most famous book on R, it covers the tidyverse and data analysis in R. Available also on the web at <https://r4ds.had.co.nz/>, including a Spanish translation, <https://es.r4ds.hadley.nz/>.
- *R Graphics Cookbook*, Chang (2018): A *cookbook* style book that covers step-by-step guides on almost all types of visualizations you may think of using ggplot2. Note: it does not stress on grammar of graphics.
- *The Effect*, Huntington-Klein (2022): This textbook, while focused on causality, describes much of the modern R development environment for statistics and econometrics. It also contains Stata and Python code.
- *Using R for Introductory Econometrics* Heiss (2020): This book is a great resource for learning R for econometrics, but it will be most useful for the Intro to Stats module.
- *Happy Git and GitHub for the useR*: This book is a great resource for learning Git and GitHub with a focus on its use for R programming.
- *R Markdown: The Definitive Guide*, Xie, Allaire, and Grolemund (2018): This book is a comprehensive guide to R Markdown, and is a great resource for learning how to write reports, papers, and presentations with R Markdown.
- *Dynamic Documents with R and knitr*, Xie (2017): The ultimate guide to knitr, the package that powers R Markdown. It allows you to understand how to fully leverage the power of R for reproducible reports.
- *Quarto website*: The Quarto website is a great resource for learning about it, which is the next generation of R Markdown.
- *The Big Book of R*: Contains links to resources on R for all topics and levels of expertise. If you want to learn something with the language that I don’t cover, this is a great place to start.

5 Software

This course requires the installation of R and RStudio. Both are free and open-source software, and are available for Windows, macOS, and Linux. You will need to install these **before** the first lecture.

- **R:** Download the latest version of R from the Comprehensive R Archive Network (CRAN) at <https://cran.r-project.org/>. Follow the instructions for your operating system.
- **RStudio:** Download the latest version of RStudio Desktop from [here](#). Follow the instructions for your operating system.

If you are having trouble installing, please post your query in the Slack channel.

You will also need Git to be able to access the course's GitHub repository. Follow the instructions from the Git & GitHub short module to make sure this software is configured correctly in your computer. As mentioned, a Git client is recommended to keep your local repository up-to-date and easily accessible to not miss any updates. For Windows users, we may use Git Bash sometimes in the course, so make sure you are familiar with the command line. Mac or Linux users can use the native terminal.

6 Communication

All communications to the instructor should be made through the course's Slack channel. We hope to respond to questions within 72 hours, but please be patient if we take longer. I do not monitor email regularly, so please use Slack for all communications if you need a timely response.

References

- Chang, Winston. 2018. *R Graphics Cookbook: Practical Recipes for Visualizing Data*. O'Reilly Media. <https://r-graphics.org/>.
- Heiss, Florian. 2020. *Using R for Introductory Econometrics*. Düsseldorf: Independently published. <https://www.urfie.net/>.
- Huntington-Klein, Nick. 2022. *The Effect: An Introduction to Research Design and Causality*. 1st edition. Boca Raton: Chapman and Hall/CRC. <https://theeffectbook.net/>.
- Wickham, Hadley, Mine Çetinkaya-Rundel, and Garrett Golemund. 2023. *R for Data Science*. "O'Reilly Media, Inc.". <https://r4ds.hadley.nz/>.
- Xie, Yihui. 2017. *Dynamic Documents with R and Knitr*. Chapman and Hall/CRC. <https://duhi23.github.io/Analisis-de-datos/Yihue.pdf>.

Xie, Yihui, Joseph J Allaire, and Garrett Grolemond. 2018. *R Markdown: The Definitive Guide*. Chapman and Hall/CRC. <https://bookdown.org/yihui/rmarkdown/>.