Robert A. Muenchen · Joseph M. Hilbe

# R for Stata Users

# Statistics and Computing

*Series Editors*
J. Chambers
D. Hand
W. Härdle

Robert A. Muenchen · Joseph M. Hilbe

# R for Stata Users

Robert A. Muenchen
University of Tennessee
Office of Information Technology
Statistical Consulting Center
916 Volunteer Blvd.
Knoxville TN 37996-0520
Stokeley Management Center
USA
muenchen.bob@gmail.com

Joseph M. Hilbe
7242 W. Heritage Way
Florence Arizona 85132
USA
hilbe@asu.edu

# Preface

While R and Stata have many features in common, their languages are quite different. Our goal in writing this book is to help you translate what you know about Stata into a working knowledge of R as quickly and easily as possible. We point out how they differ using terminology with which you are familiar and we include many Stata terms in the table of contents and index. You can find any R function by looking up its counterpart in Stata and vice versa. We provide many example programs done in R and Stata so that you can see how they compare topic by topic.

When finished, you should be able to use R to:

- Read data from various types of text files and Stata data sets.
- Manage your data through transformations, recodes, and combining data sets from both the add-cases and add-variables approaches and restructuring data from wide to long formats and vice versa.
- Create publication quality graphs including bar, histogram, pie, line, scatter, regression, box, error bar, and interaction plots.
- Perform the basic types of analyses to measure strength of association and group differences and be able to know where to turn to cover much more complex methods.

## Who This Book Is For

This book is, of course, for people who already know Stata. It may also be useful to R users wishing to learn Stata. However, we explain none of the Stata programs, only the R ones and how the packages differ, so it is not ideal for that purpose.

This book is based on *R for SAS and SPSS Users* [34]. However, there is quite a bit of additional material covered here, and, of course, the comparative coverage is completely different.

## Who This Book Is Not For

We make no effort to teach statistics or graphics. Although we briefly state the goal and assumptions of each analysis, we do not cover their formulas or derivations. We have more than enough to discuss without tackling those topics too. This is also not a book about writing R functions, it is about using the thousands that already exist. We will write only a few very short functions. If you want to learn more about writing functions, we recommend John Chamber's *Software for Data Analysis: Programming with R* [5]. However, if you know Stata, reading this book should ease your transition to more complex books like that.

## Practice Data Sets and Programs

All of the programs, data sets, and files that we use in this book are available for download at `http://r4stats.com`. A file containing corrections and clarifications is also available there.

## Acknowledgments

We are very grateful for the many people who have helped make this book possible, including the developers of the S language on which R is based, Rick Becker, John Chambers, and Allan Wilks; the people who started R itself, Ross Ihaka and Robert Gentleman; the many other R developers for providing such wonderful tools for free and all the R-help participants who have kindly answered so many questions. Virtually all of the examples we present here are modestly tweaked versions of countless posts to the R-help discussion list, as well as a few Statalist posts. All we add is the selection, organization, explanation, and comparison.

    We are especially grateful to the people who provided advice, caught typos, and suggested improvements, including Raymond R. Balise, Patrick Burns, Peter Flom, Chun Huang, Martin Gregory, Warren Lambert, Mathew Marler, Ralph O'Brien, Wayne Richter, Charilaos Skiadas, Andreas Stefik, Phil Spector, Michael Wexler, Graham Williams, Andrew Yee, and several anonymous reviewers.

    A special thanks goes to Hadley Wickham, who provided much guidance on his `ggplot2` graphics package. Thanks to Gabor Grothendieck, Lauri Nikkinen, and Marc Schwarz and for the R-Help help discussion that led to Section 10.14: "Selecting First or Last Observations per Group." Thanks to Gabor Grothendieck also for a detailed discussion that lead to Section 10.4: "Multiple Conditional Transformations." Thanks to Michael A. McGuire for his assistance with all things Macintosh.

The first author is grateful to his wife, Carla Foust, and sons Alexander and Conor, who put up with many lost weekends as he wrote this book.

The second author wishes to thank Springer editor John Kimmel for suggesting his participation in this project and his wife, Cheryl, children Heather, Michael and Mitchell, and Sirr for their patience while he spent time away from them working on this book.

*Robert A. Muenchen*
muenchen.bob@gmail.com
Knoxville, Tennessee
January 2010

*Joseph M. Hilbe*
hilbe@asu.edu
Florence, Arizona
January 2010

## About the Authors

Robert A. Muenchen is a consulting statistician and author of the book, *R for SAS and SPSS Users* [34]. He is currently the manager of Research Computing Support (formerly the Statistical Consulting Center) at the University of Tennessee. Bob has conducted research for a variety of public and private organizations and has co-authored over 50 articles in scientific journals and conference proceedings.

Bob has served on the advisory boards of the SAS Institute, SPSS Inc., the Statistical Graphics Corporation, and *PC Week Magazine*. His suggested improvements have been incorporated into SAS, SPSS, JMP, STATGRAPHICS, and several R packages.

His research interests include statistical computing, data graphics and visualization, text analysis, data mining, psychometrics, and resampling.

Joseph M. Hilbe is Solar System Ambassador with NASA/Jet Propulsion Laboratory, California Institute of Technology, an adjunct professor of statistics at Arizona State, and emeritus professor at the University of Hawaii. He is an elected Fellow of the American Statistical Association and of the Royal Statistical Society and is an elected member of the International Statistical Institute.

Professor Hilbe was the first editor of the *Stata Technical Bulletin*, later to become the *Stata Journal*, and was one of Stata Corporation's first senior statisticians (1991–1993). Hilbe is also the author of a number of textbooks,

including *Logistic Regression Models* [21], *Negative Binomial Regression* [23], and with J. Hardin, *Generalized Linear Models and Extensions*, 2nd ed. [18] and *Generalized Estimating Equations* [19].

# Contents

# List of Tables

# List of Figures

# 1

# Introduction

## 1.1 Overview

R [38] is a powerful and flexible environment for research computing. Written
by Ross Ihaka, Robert Gentleman (hence the name "R"), the R Core Develop-
ment Team, and an army of volunteers, R provides a wider range of analytical
and graphical commands than any other software. The fact that this level of
power is available free of charge has dramatically changed the landscape of
research software.

R is a variation of the S language, developed by John Chambers, Rick
Becker, and others at Bell Labs[1]. The Association of Computing Machinery
presented John Chambers with a Software System Award and said that the S
language "...*will forever alter the way people analyze, visualize, and manip-
ulate data...*" and went on to say that it is "...*an elegant, widely accepted,
and enduring software system, with conceptual integrity....*" The original S
language is still commercially available as Tibco Spotfire S+. Most programs
written in the S language will run in R.

Stata, a product of Stata Corporation, has not yet incorporated an inter-
face to R in its software, but users have already posted programs to use R
within the Stata environment. It is expected that more facilities of this sort
will be developed in the near future.

For each aspect of R we discuss, we will compare and contrast it with Stata.
Many of the topics end with example programs that do almost identical things
in both software applications. R programs are often longer than similar Stata
code, but this is typically the case because R functions are more specific than
Stata commands.

Many R functions will appear familiar to Stata users; that is, R functions
such as `lm` or `glm` will appear somewhat similar to Stata's `regress` and `glm`
commands. There are other aspects of the two languages, however, that may

---

[1] For a fascinating history of S and R, see Appendix A of *Software for Data
Analysis: Programming with R* [5].

appear more confusing at first. We hope to ease that confusion by focusing on both the similarities and differences between R and Stata in this text. When we examine a particular analysis (e.g., comparing two groups with a t-test) someone who knows Stata will have very little trouble figuring out what R is doing. However, the basics of the R language are very different, so that is where we will spend the majority of our time.

We introduce topics in a carefully chosen order, so it is best to read from beginning to end the first time through, even if you think you do not need to know a particular topic. Later you can skip directly to the section you need. We include a fair amount of redundancy on key topics to help teach those topics and to make it easier to read just one section as a future reference. The glossary in Appendix A defines R concepts in terms that Stata users will understand and provides parallel definitions using R terminology.

## 1.2 Similarities Between R and Stata

Stata is an excellent statistics package. One of the authors has used Stata for over 20 years and has authored many Stata commands.

Perhaps more than any other two research computing environments, R and Stata share many of the features that make them outstanding:

- Both include rich programming languages designed for writing new analytic methods, not just a set of prewritten commands.
- Both contain extensive sets of analytic commands written in their own languages.
- The pre-written commands in R, and most in Stata, are visible and open for you to change as you please.
- Both save command or function output in a form you can easily use as input to further analysis.
- Both do modeling in a way that allows you to readily apply your models for tasks such as making predictions on new data sets. Stata calls these *postestimation commands* and R calls them *extractor functions*.
- In both, when you write a new command, it is on an equal footing with commands written by the developers. There are no additional "Developer's Kits" to purchase.
- Both have legions of devoted users who have written numerous extensions and who continue to add the latest methods many years before their competitors.
- Both can search the Internet for user-written commands and download them automatically to extend their capabilities quickly and easily.
- Both hold their data in the computer's main memory, offering speed but limiting the amount of data they can handle.

# 1.3 Why Learn R?

With so many similarities, if you already know Stata, why should you bother to learn R?

- To augment Stata; i.e. to be able to perform statistical analyses that are not available in Stata, but which are available in R. R offers a *vast* number of analytical methods. There are now over 3,000 add-on packages available for R and this number is growing at an exponential rate. Therefore, knowing both gives you a much greater range of tools for analyzing data.
- To stay current with new analytic methods. The majority of statistics textbooks, and journal articles, now being published use either Stata or R for examples. R appears to be used more in many journals. Stata users not understanding R are therefore not able to learn as much from texts or articles using R for examples than they would be if they understood the language.
- If you continue to do all of your data management in Stata, you can learn just enough R to import your data and run the procedures you need.
- R is directly accessible from inside many statistics packages. SAS, SPSS, and STATISTICA offer the ability to run R programs from within their software. This means that when developers write programs in R, they are assured a very wide audience. Roger Newson has written an interface [36] between Stata and R that provides some of this ability. We expect to see more done on this topic in the near future.
- R has been object-oriented since its first version. Many of its commands sense the types of data structures you have and do the best thing for each. For example, once you tell it that gender is a categorical variable, it will take statistically proper actions if you use it as a linear regression predictor. At the time of publication, Stata Corporation had just announced its future move toward object orientation.
- Both languages consist of a core set of functions that are written in the C language. However, only developers at Stata Corporation can modify its most fundamental commands. Every aspect of R is open for anyone to modify in any way they like. This complete flexibility attracts many developers.
- Both R and Stata offer graphics that are flexible, easy to use, and of high quality. However, R also offers the very flexible and powerful Grammar of Graphics approach. As we will see, developers have even gone so far as replacing R's core graphical system.
- R is free. This means, of course, that you can use it for free, but it also means developers know that their work is available to everyone. That helps attract developers and is a major reason that there are so many add-on packages for it.

## 1.4 Is R Accurate?

When people first learn of R, one of their first questions is "Can a package written by volunteers be as accurate as one written by a large corporation?" People envision a lone programmer competing against a large corporate team. Having worked closely with several software companies over the years, we can assure you that this is not the case. A particular procedure is usually written by one programmer, even at Stata Corporation. A thorough testing process is then carried out by a few people within the company and then more thoroughly by Stata users on publication of the new command or function.

The R Development Core Team runs each release of R through validation suites that have known correct answers to ensure accurate results. They also go through "Alpha," "Beta," and "Release Candidate" testing phases, which are open to the public. Each phase has tighter restrictions on modifications of R. Finally, the production version is released. The details of this process are provided in R: Regulatory Compliance and Validation Issues, A Guidance Document for the Use of R in Regulated Clinical Trial Environments, available at `http://www.r-project.org/doc/R-FDA.pdf` [11].

When bugs are found in Stata, the developers typically make a fix within days. Users are in continual communication with other users and developers through the Statalist. An average of 100 communications are posted daily. Questions are answered by other users or by Stata staff.

R also has open discussions of its known bugs and R's developers fix them quickly too. However, software of the complexity of Stata and R will never be completely free of errors, regardless of its source.

## 1.5 What About Tech Support?

If a package is free, who supports it?

Stata users may call toll-free or e-mail technical support for problems they experience with the software or for advice on how to run various software commands. The response is near immediate, with a day delay in response being on the high side. Even experienced Stata users sometimes require technical advice for new commands or functions or have difficulties learning new areas of statistics or new methodologies (e.g. the matrix programming). We have always found support to he helpful and friendly.

You can also get support through the Stata Listserver, where it is normal to get assistance from someone the very day you post your request.

R's main source of support is the R-help mailing list. Other users and often developers themselves will often provide immediate help. Sometimes you may obtain different answers from various responders, but that is part of the nature of statistics. For details on the various R e-mail support lists, see Chapter 4, "Help and Documentation."

There are several commercial versions of R available, and the companies that sell them do provide phone support. Here are some of these companies and their web sites:

XL-Solutions Corp., `http://www.experience-rplus.com/`

Revolution Computing, Inc., `http://www.revolution-computing.com/`

Random Technologies, LLC, `http://random-technologies-llc.com/`

## 1.6 Getting Started Quickly

If you wish to start using R quickly, you can do so by reading fewer than 50 pages of this book. Since you have Stata to do your basic descriptive statistics, you are likely to need R's modeling functions. Here are the steps you can follow to use them.

1. Read the remainder of this chapter and Chapter 2, "Installing and Updating R." Download and install R on your computer.
2. Read the part of Chapter 3, "Running R," that covers your operating system.
3. In Chapter 5, "Programming Language Basics," read Section 5.3.2 about factors, and Section 5.3.3 about data frames.
4. Also in Chapter 5, read Section 5.6.1, "Controlling Functions with Arguments," and Section 5.6.2, "Controlling Functions with Formulas," including Table 5.1, "Example formulas in Stata and R."
5. Read Section 6.6, "Importing Data from Stata."

After reading the pages above, do all your data management in Stata, stripping out observations containing any missing values. Then write out only the variables and observations you need to a comma separated values file, mydata.csv. Assuming your variables are named y, x1, x2,..., your entire R program will look something like this:

```
library("Hmisc")  # Contains stata.get function.
library("OtherLibrariesYouNeed")  # If you need any.
mydata <- stata.get("mydata.dta") # imports your Stata file
mymodel <- TheFunctionYouNeed( y ~ x1+x2, data=mydata )
summary(mymodel)
plot(mymodel) # if your function does plots.
```

## 1.7 Programming Conventions

Although R has many ways to generate practice data and has a variety of example data sets, we will use a tiny practice data set that is easy to enter. We can then manipulate and print it repeatedly so that you can clearly see the changes.