# Advanced R

# Chapman & Hall/CRC
# The R Series

## Series Editors

**John M. Chambers**
Department of Statistics
Stanford University
Stanford, California, USA

**Torsten Hothorn**
Division of Biostatistics
University of Zurich
Switzerland

**Duncan Temple Lang**
Department of Statistics
University of California, Davis
Davis, California, USA

**Hadley Wickham**
RStudio
Boston, Massachusetts, USA

## Aims and Scope

This book series reflects the recent rapid growth in the development and application of R, the programming language and software environment for statistical computing and graphics. R is now widely used in academic research, education, and industry. It is constantly growing, with new versions of the core software released regularly and more than 5,000 packages available. It is difficult for the documentation to keep pace with the expansion of the software, and this vital book series provides a forum for the publication of books covering many aspects of the development and application of R.

The scope of the series is wide, covering three main threads:
- Applications of R to specific disciplines such as biology, epidemiology, genetics, engineering, finance, and the social sciences.
- Using R for the study of topics of statistical methodology, such as linear and mixed modeling, time series, Bayesian methods, and missing data.
- The development of R, including programming, building packages, and graphics.

The books will appeal to programmers and developers of R software, as well as applied statisticians and data analysts in many fields. The books will feature detailed worked examples and R code fully integrated into the text, ensuring their usefulness to researchers, practitioners and students.

# Published Titles

**Stated Preference Methods Using R**, *Hideo Aizaki, Tomoaki Nakatani, and Kazuo Sato*

**Using R for Numerical Analysis in Science and Engineering**, *Victor A. Bloomfield*

**Event History Analysis with R**, *Göran Broström*

**Computational Actuarial Science with R**, *Arthur Charpentier*

**Statistical Computing in C++ and R**, *Randall L. Eubank and Ana Kupresanin*

**Reproducible Research with R and RStudio**, *Christopher Gandrud*

**Introduction to Scientific Programming and Simulation Using R, Second Edition**, *Owen Jones, Robert Maillardet, and Andrew Robinson*

**Nonparametric Statistical Methods Using R**, *John Kloke and Joseph McKean*

**Displaying Time Series, Spatial, and Space-Time Data with R**, *Oscar Perpiñán Lamigueiro*

**Programming Graphical User Interfaces with R**, *Michael F. Lawrence and John Verzani*

**Analyzing Sensory Data with R**, *Sébastien Lê and Theirry Worch*

**Analyzing Baseball Data with R**, *Max Marchi and Jim Albert*

**Growth Curve Analysis and Visualization Using R**, *Daniel Mirman*

**R Graphics, Second Edition**, *Paul Murrell*

**Multiple Factor Analysis by Example Using R**, *Jérôme Pagès*

**Customer and Business Analytics: Applied Data Mining for Business Decision Making Using R**, *Daniel S. Putler and Robert E. Krider*

**Implementing Reproducible Research**, *Victoria Stodden, Friedrich Leisch, and Roger D. Peng*

**Using R for Introductory Statistics, Second Edition**, *John Verzani*

**Advanced R**, *Hadley Wickham*

**Dynamic Documents with R and knitr**, *Yihui Xie*

# Advanced R

## Hadley Wickham

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

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

*To Jeff, who makes me happy, and who made sure I had a life outside this book.*

# *Contents*

**9  Debugging, condition handling, and defensive programming**                                          **149**

## III    Computing on the language     257

## 13 Non-standard evaluation     259

## 14 Expressions     281