# Use R!

# Use R!

Dianne Cook   Deborah F. Swayne

# Interactive and Dynamic Graphics for Data Analysis

## With R and GGobi

With Contributions by Andreas Buja, Duncan Temple Lang,
Heike Hofmann, Hadley Wickham, and Michael Lawrence

Springer

Dianne Cook
Department of Statistics
Iowa State University
325 Snedecor Hall
Ames, IA 50011-1210
dicook@iastate.edu

Deborah F. Swayne
AT & T Labs - Research
Shannon Laboratory
180 Park Avenue
Florham Park, NJ 07932-1049
dfs@research.att.com

*Series Editors:*
Robert Gentleman
Program in Computational Biology
Division of Public Health Sciences
Fred Hutchinson Cancer Research Center
1100 Fairview Ave. N, M2-B876
Seattle, Washington 981029-1024
USA

Kurt Hornik
Department für Statistik und Mathematik
Wirtschaftsuniversität Wien Augasse 2-6
A-1090 Wien
Austria

Giovanni Parmigiani
The Sidney Kimmel Comprehensive Cancer Center at Johns Hopkins University
550 North Broadway
Baltimore, MD 21205-2011
USA

*Cover illustration*: Brushing a cluster that was found by using the grand tour to explore two-dimensional
projections of five-dimensional data on Australian crabs.

Printed on acid-free paper.

# Preface

This book is about using interactive and dynamic plots on a computer screen as part of data exploration and modeling, both alone and as a partner with static graphics and non-graphical computational methods. The area of interactive and dynamic data visualization emerged within statistics as part of research on exploratory data analysis in the late 1960s, and it remains an active subject of research today, as its use in practice continues to grow. It now makes substantial contributions within computer science as well, as part of the growing fields of information visualization and data mining, especially visual data mining.

The material in this book includes:

- An introduction to data visualization, explaining how it differs from other types of visualization.
- A description of our toolbox of interactive and dynamic graphical methods.
- An approach for exploring missing values in data.
- An explanation of the use of these tools in cluster analysis and supervised classification.
- An overview of additional material available on the web.
- A description of the data used in the analyses and exercises.

The book's examples use the software R and GGobi. R (Ihaka & Gentleman 1996, R Development Core Team 2006) is a free software environment for statistical computing and graphics; it is most often used from the command line, provides a wide variety of statistical methods, and includes high–quality static graphics. R arose in the Statistics Department of the University of Auckland and is now developed and maintained by a global collaborative effort. It began as a re-implementation of the S language and statistical computing environment (Becker & Chambers 1984) first developed at Bell Laboratories before the breakup of AT&T.

GGobi (Swayne, Temple Lang, Buja & Cook 2003) is free software for interactive and dynamic graphics; it can be operated using a command-line interface or from a graphical user interface (GUI). When GGobi is used as a

stand-alone tool, only the GUI is used; when it is used with R, via the `rggobi`
(Temple Lang, Swayne, Wickham & Lawrence 2006) package, a command-line
interface is used along with the GUI. GGobi is a descendant of two earlier pro-
grams: XGobi (Swayne, Cook & Buja 1992, Buja, Cook & Swayne 1996) and,
before that, Dataviewer (Buja, Hurley & McDonald 1986, Hurley 1987). Many
of the examples that follow might be reproduced with other software such
as S-PLUS®, SAS JMP®, DataDesk®, Mondrian, MANET, and Spotfire®.
However, GGobi is unique because it offers tours (rotations of data in higher
than 3D), complex linking between plots using categorical variables, and the
tight connection with R.

*Web resources*

The web site which accompanies the book contains sample datasets and
R code, movies demonstrating the interactive and dynamic graphic meth-
ods, and additional chapters. It can be reached through the GGobi web site:

<div align="center">

`http://www.ggobi.org`

</div>

The R software is available from:

<div align="center">

`http://www.R-project.org`

</div>

Both web sites include source code as well as binaries for various operating
systems (Linux®, Windows®, Mac OS X®) and allow users to sign up for
mailing lists and browse mailing list archives. The R web site offers a wealth
of documentation, including an introduction to R and a partially annotated
list of books offering more instruction. [Widely read books include Dalgaard
(2002), Venables & Ripley (2002), and Maindonald & Braun (2003).] The
GGobi web site includes an introductory tutorial, a list of papers, and several
movies.

*How to use this book*

The language in the book is aimed at later year undergraduates, beginning
graduate students, and graduate students in any discipline needing to analyze
their own multivariate data. It is suitable reading for an industry statisti-
cian, engineer, bioinformaticist, or computer scientist with some knowledge
of basic data analysis and a need to analyze high-dimensional data. It also
may be useful for a mathematician who wants to visualize high-dimensional
structures.

The end of each chapter contains exercises to help practice the methods
discussed in the chapter. The book may be used as a text in a class on statis-
tical graphics, exploratory data analysis, visual data mining, or information
visualization. It might also be used as an adjunct text in a course on multi-
variate data analysis or data mining.

This book has been closely tied to a particular software implementation
so that you can actively use the methods as you read about them, to learn
and experiment with interactive and dynamic graphics. The plots and writ-
ten explanations in the book are no substitute for personal experience. We
strongly urge the reader to go through this book sitting near a computer with

GGobi, R, and `rggobi` installed, following along with the examples. If you do not wish to install the software, then the next best choice is to watch the accompanying movies demonstrating the examples in the text.

If you have not used GGobi before, then visit the web site, watch the movies, download the manual, and work through the tutorial; the same advice applies for those unfamiliar with R: Visit the R web site and learn the basics.

As you read the book, try things out for yourself. Take your time, and have fun!

*Acknowledgments*

Dianne Cook                    Deborah F. Swayne
Iowa State University          AT&T Labs – Research

July 2007

# Contents

# Technical Notes

*R code*

The R code in this book, denoted by `typewriter font`, and the more extensive code on the web site, has been tested on version 2.4.0 of R, version 2.1.5 of GGobi, and version 2.1.5 of `rggobi`. Updates will be available on the web site as they are needed.

*Figures*

The figures in this book were produced in a variety of ways, and the files and code to reproduce them are all available on the book's web site. Some were produced directly in R. Some were produced using both GGobi and R, and the process of converting GGobi views into publication graphics deserves an explanation.

When we arrive at a GGobi view we want to include in a paper or book, we use the `Save Display Description` item on GGobi's `Tools` menu to generate a file containing an S language description of the display. We read the file into R using the R package `DescribeDisplay` (Wickham 2006*b*), like this:

```
> library(DescribeDisplay)
> d <- dd_load("fig.R")
```

We create the publication-quality graphic using either that package's plot method or another R package, `ggplot` (Wickham 2006*c*), like this:

```
> plot(d)
```

or

```
> p <- ggplot(d)
> print(p)
```

Figure 0.1 illustrates the differences with a trio of representations of the same bivariate scatterplot. The picture at left is a screen dump of a GGobi display. Such images are not usually satisfactory for publication for several

**Fig. 0.1.** Sample plots produced from GGobi in different ways: **(left)** a simple screen dump; **(middle)** a plot produced using the plot method of the R package `DescribeDisplay`; **(right)** a plot made using the R package `ggplot`.

reasons, the most obvious of which is the lack of resolution. The second picture was produced using `DescribeDisplay`'s plot method, which reproduces the plotting region of the view with pretty good fidelity. We used this method to produce most of the one–dimensional and two-dimensional tour pictures in this book. The third picture was produced using `ggplot`, which adds axis ticks, labels and grid lines. We used it to produce nearly all the bivariate scatterplots of GGobi views in this book.

# List of Figures