

Introduction to Probability and Statistics Using R

Second Edition

G. Jay Kerns

October 3, 2011

IP_SUR: Introduction to Probability and Statistics Using R
Copyright © 2011 G. Jay Kerns ISBN: 978-0-557-24979-4

Permission is granted to copy, distribute and/or modify this document under the terms of the GNU Free Documentation License, Version 1.3 or any later version published by the Free Software Foundation; with no Invariant Sections, no Front-Cover Texts, and no Back-Cover Texts. A copy of the license is included in the section entitled “GNU Free Documentation License”.

Date: October 3, 2011

Contents

Preface to the Second Edition	v
Preface to the First Edition	vii
List of Figures	xiii
List of Tables	xv
1 An Introduction to R	xvii
1.1 Downloading and Installing R	xvii
1.2 Communicating with R	xix
1.3 Basic R Operations and Concepts	xxi
1.4 Getting Help	xxvii
1.5 External Resources	xxix
1.6 Other Tips	xxix
1.7 Exercises	xxxi

Preface to the Second Edition

Preface to the First Edition

This book was expanded from lecture materials I use in a one semester upper-division undergraduate course entitled *Probability and Statistics* at Youngstown State University. Those lecture materials, in turn, were based on notes that I transcribed as a graduate student at Bowling Green State University. The course for which the materials were written is 50-50 Probability and Statistics, and the attendees include mathematics, engineering, and computer science majors (among others). The catalog prerequisites for the course are a full year of calculus.

The book can be subdivided into three basic parts. The first part includes the introductions and elementary *descriptive statistics*; I want the students to be knee-deep in data right out of the gate. The second part is the study of *probability*, which begins at the basics of sets and the equally likely model, journeys past discrete/continuous random variables, and continues through to multivariate distributions. The chapter on sampling distributions paves the way to the third part, which is *inferential statistics*. This last part includes point and interval estimation, hypothesis testing, and finishes with introductions to selected topics in applied statistics.

I usually only have time in one semester to cover a small subset of this book. I cover the material in Chapter 2 in a class period that is supplemented by a take-home assignment for the students. I spend a lot of time on Data Description, Probability, Discrete, and Continuous Distributions. I mention selected facts from Multivariate Distributions in passing, and discuss the meaty parts of Sampling Distributions before moving right along to Estimation (which is another chapter I dwell on considerably). Hypothesis Testing goes faster after all of the previous work, and by that time the end of the semester is in sight. I normally choose one or two final chapters (sometimes three) from the remaining to survey, and regret at the end that I did not have the chance to cover more.

In an attempt to be correct I have included material in this book which I would normally not mention during the course of a standard lecture. For instance, I normally do not highlight the intricacies of measure theory or integrability conditions when speaking to the class. Moreover, I often stray from the matrix approach to multiple linear regression because many of my students have not yet been formally trained in linear algebra. That being said, it is important to me for the students to hold something in their hands which acknowledges the world of mathematics and statistics beyond the classroom, and which may be useful to them for many semesters to come. It also mirrors my own experience as a student.

The vision for this document is a more or less self contained, essentially complete,

correct, introductory textbook. There should be plenty of exercises for the student, with full solutions for some, and no solutions for others (so that the instructor may assign them for grading). By Sweave's dynamic nature it is possible to write randomly generated exercises and I had planned to implement this idea already throughout the book. Alas, there are only 24 hours in a day. Look for more in future editions.

Seasoned readers will be able to detect my origins: *Probability and Statistical Inference* by Hogg and Tanis [?], *Statistical Inference* by Casella and Berger [?], and *Theory of Point Estimation* and *Testing Statistical Hypotheses* by Lehmann [?, ?]. I highly recommend each of those books to every reader of this one. Some R books with “introductory” in the title that I recommend are *Introductory Statistics with R* by Dalgaard [?] and *Using R for Introductory Statistics* by Verzani [?]. Surely there are many, many other good introductory books about R, but frankly, I have tried to steer clear of them for the past year or so to avoid any undue influence on my own writing.

I would like to make special mention of two other books: *Introduction to Statistical Thought* by Michael Lavine [?] and *Introduction to Probability* by Grinstead and Snell [?]. Both of these books are *free* and are what ultimately convinced me to release *IP&UR* under a free license, too.

Please bear in mind that the title of this book is “Introduction to Probability and Statistics Using R”, and not “Introduction to R Using Probability and Statistics”, nor even “Introduction to Probability and Statistics and R Using Words”. The people at the party are Probability and Statistics; the handshake is R. There are several important topics about R which some individuals will feel are underdeveloped, glossed over, or wantonly omitted. Some will feel the same way about the probabilistic and/or statistical content. Still others will just want to learn R and skip all of the mathematics.

Despite any misgivings: here it is, warts and all. I humbly invite said individuals to take this book, with the GNU Free Documentation License (GNU-FDL) in hand, and make it better. In that spirit there are at least a few ways in my view in which this book could be improved.

Better data. The data analyzed in this book are almost entirely from the `datasets` package in base R, and here is why:

- I made a conscious effort to minimize dependence on contributed packages,
- The data are instantly available, already in the correct format, so we need not take time to manage them, and
- The data are *real*.

I made no attempt to choose data sets that would be interesting to the students; rather, data were chosen for their potential to convey a statistical point. Many of the data sets are decades old or more (for instance, the data used to introduce simple linear regression are the speeds and stopping distances of cars in the 1920's).

In a perfect world with infinite time I would research and contribute recent, *real* data in a context crafted to engage the students in *every* example. One day I hope to stumble over said time. In the meantime, I will add new data sets incrementally as time permits.

More proofs. I would like to include more proofs for the sake of completeness (I understand that some people would not consider more proofs to be improvement). Many proofs have been skipped entirely, and I am not aware of any rhyme or reason to the current omissions. I will add more when I get a chance.

More and better~graphics. I have not used the `ggplot2` package [?] because I do not know how to use it yet. It is on my to-do list.

More and better exercises. There are only a few exercises in the first edition simply because I have not had time to write more. I have toyed with the `exams` package [?] and I believe that it is a right way to move forward. As I learn more about what the package can do I would like to incorporate it into later editions of this book.

About This Document

IPSUR contains many interrelated parts: the *Document*, the *Program*, the *Package*, and the *Ancillaries*. In short, the *Document* is what you are reading right now. The *Program* provides an efficient means to modify the Document. The *Package* is an R package that houses the Program and the Document. Finally, the *Ancillaries* are extra materials that reside in the Package and were produced by the Program to supplement use of the Document. We briefly describe each of them in turn.

The Document

The *Document* is that which you are reading right now – IPSUR’s *raison d’être*. There are transparent copies (nonproprietary text files) and opaque copies (everything else). See the GNU-FDL in Appendix ?? for more precise language and details.

IPSUR.tex is a transparent copy of the Document to be typeset with a L^AT_EX distribution such as MikTeX or T_EX Live. Any reader is free to modify the Document and release the modified version in accordance with the provisions of the GNU-FDL. Note that this file cannot be used to generate a randomized copy of the Document. Indeed, in its released form it is only capable of typesetting the exact version of IPSUR which you are currently reading. Furthermore, the `.tex` file is unable to generate any of the ancillary materials.

IPSUR-xxx.eps, IPSUR-xxx.pdf are the image files for every graph in the Document. These are needed when typesetting with L^AT_EX.

Contents

IPSUR.pdf is an opaque copy of the Document. This is the file that instructors would likely want to distribute to students.

IPSUR.dvi is another opaque copy of the Document in a different file format.

The Program

The *Program* includes `IPSUR.lyx` and its nephew `IPSUR.Rnw`; the purpose of each is to give individuals a way to quickly customize the Document for their particular purpose(s).

IPSUR.lyx is the source LyX file for the Program, released under the GNU General Public License (GNU GPL) Version 3. This file is opened, modified, and compiled with LyX, a sophisticated open-source document processor, and may be used (together with Sweave) to generate a randomized, modified copy of the Document with brand new data sets for some of the exercises and the solution manuals (in the Second Edition). Additionally, LyX can easily activate/deactivate entire blocks of the document, *e.g.* the *proofs* of the theorems, the student *solutions* to the exercises, or the instructor *answers* to the problems, so that the new author may choose which sections (s)he would like to include in the final Document (again, Second Edition). The `IPSUR.lyx` file is all that a person needs (in addition to a properly configured system – see Appendix ??) to generate/compile/export to all of the other formats described above and below, which includes the ancillary materials `IPSUR.Rdata` and `IPSUR.R`.

IPSUR.Rnw is another form of the source code for the Program, also released under the GNU GPL Version 3. It was produced by exporting `IPSUR.lyx` into R/Sweave format (`.Rnw`). This file may be processed with Sweave to generate a randomized copy of `IPSUR.tex` – a transparent copy of the Document – together with the ancillary materials `IPSUR.Rdata` and `IPSUR.R`. Please note, however, that `IPSUR.Rnw` is just a simple text file which does not support many of the extra features that LyX offers such as WYSIWYM editing, instantly (de)activating branches of the manuscript, and more.

The Package

There is a contributed package on CRAN, called `IPSUR`. The package affords many advantages, one being that it houses the Document in an easy-to-access medium. Indeed, a student can have the Document at his/her fingertips with only three commands:

Another advantage goes hand in hand with the Program's license; since `IPSUR` is free, the source code must be freely available to anyone that wants it. A package hosted on CRAN allows the author to obey the license by default.

A much more important advantage is that the excellent facilities at R-Forge are building and checking the package daily against patched and development versions of the absolute latest pre-release of R. If any problems surface then I will know about it within 24 hours.

And finally, suppose there is some sort of problem. The package structure makes it *incredibly* easy for me to distribute bug-fixes and corrected typographical errors. As an author I can make my corrections, upload them to the repository at R-Forge, and they will be reflected *worldwide* within hours. We aren't in Kansas anymore, Toto.

Ancillary Materials

These are extra materials that accompany `IPSUR`. They reside in the `/etc` subdirectory of the package source.

IPSUR.RData is a saved image of the R workspace at the completion of the Sweave processing of `IPSUR`. It can be loaded into memory with `File ▶ Load Workspace` or with the command `load("/path/to/IPSUR.Rdata")`. Either method will make every single object in the file immediately available and in memory. In particular, the data `BLANK` from Exercise `BLANK` in Chapter `BLANK` on page `BLANK` will be loaded. Type `BLANK` at the command line (after loading `IPSUR.RData`) to see for yourself.

IPSUR.R is the exported R code from `IPSUR.Rnw`. With this script, literally every R command from the entirety of `IPSUR` can be resubmitted at the command line.

Notation

We use the notation `x` or `stem.leaf` notation to denote objects, functions, *etc.*. The sequence `Statistics ▶ Summaries ▶ Active Dataset` means to click the `Statistics` menu item, next click the `Summaries` submenu item, and finally click `Active Dataset`.

Acknowledgements

This book would not have been possible without the firm mathematical and statistical foundation provided by the professors at Bowling Green State University, including Drs. Gábor Székely, Craig Zirbel, Arjun K. Gupta, Hanfeng Chen, Truc Nguyen, and James Albert. I would also like to thank Drs. Neal Carothers and Kit Chan.

I would also like to thank my colleagues at Youngstown State University for their support. In particular, I would like to thank Dr. G. Andy Chang for showing me what it means to be a statistician.

I would like to thank Richard Heiberger for his insightful comments and improvements to several points and displays in the manuscript.

Contents

Finally, and most importantly, I would like to thank my wife for her patience and understanding while I worked hours, days, months, and years on a *free book*. Looking back, I can't believe I ever got away with it.

List of Figures

List of Tables

1 An Introduction to R

1.1 Downloading and Installing R

The instructions for obtaining R largely depend on the user's hardware and operating system. The R Project has written an R Installation and Administration manual with complete, precise instructions about what to do, together with all sorts of additional information. The following is just a primer to get a person started.

1.1.1 Installing R

Visit one of the links below to download the latest version of R for your operating system:

Microsoft Windows: <http://cran.r-project.org/bin/windows/base/>

MacOS: <http://cran.r-project.org/bin/macosx/>

Linux: <http://cran.r-project.org/bin/linux/>

On Microsoft Windows, click the `R-x.y.z.exe` installer to start installation. When it asks for “Customized startup options”, specify Yes. In the next window, be sure to select the SDI (single document interface) option; this is useful later when we discuss three dimensional plots with the `rgl` package [?].

Installing R on a USB drive (Windows)

With this option you can use R portably and without administrative privileges. There is an entry in the R for Windows FAQ about this. Here is the procedure I use:

1. Download the Windows installer above and start installation as usual. When it asks *where* to install, navigate to the top-level directory of the USB drive instead of the default C drive.
2. When it asks whether to modify the Windows registry, uncheck the box; we do NOT want to tamper with the registry.
3. After installation, change the name of the folder from `R-x.y.z` to just plain R. (Even quicker: do this in step 1.)

4. [Download this shortcut](#) and move it to the top-level directory of the USB drive, right beside the R folder, not inside the folder. Use the downloaded shortcut to run R.

Steps 3 and 4 are not required but save you the trouble of navigating to the `R-x.y.z/bin` directory to double-click `Rgui.exe` every time you want to run the program. It is useless to create your own shortcut to `Rgui.exe`. Windows does not allow shortcuts to have relative paths; they always have a drive letter associated with them. So if you make your own shortcut and plug your USB drive into some *other* machine that happens to assign your drive a different letter, then your shortcut will no longer be pointing to the right place.

1.1.2 Installing and Loading Add-on Packages

There are *base* packages (which come with R automatically), and *contributed* packages (which must be downloaded for installation). For example, on the version of R being used for this document the default base packages loaded at startup are

```
getOption("defaultPackages")
```

```
[1] "datasets" "utils"      "grDevices" "graphics"
[5] "stats"    "methods"
```

The base packages are maintained by a select group of volunteers, called R Core. In addition to the base packages, there are literally thousands of additional contributed packages written by individuals all over the world. These are stored worldwide on mirrors of the Comprehensive R Archive Network, or CRAN for short. Given an active Internet connection, anybody is free to download and install these packages and even inspect the source code.

To install a package named `foo`, open up R and type `install.packages("foo")`. To install `foo` and additionally install all of the other packages on which `foo` depends, instead type `install.packages("foo", depends = TRUE)`.

The general command `install.packages()` will (on most operating systems) open a window containing a huge list of available packages; simply choose one or more to install.

No matter how many packages are installed onto the system, each one must first be loaded for use with the `library` function. For instance, the `foreign` package [?] contains all sorts of functions needed to import data sets into R from other software such as SPSS, SAS, *etc.* But none of those functions will be available until the command `library(foreign)` is issued.

Type `library()` at the command prompt (described below) to see a list of all available packages in your library.

For complete, precise information regarding installation of R and add-on packages, see the [R Installation and Administration manual](#).

1.2 Communicating with R

1.2.1 One line at a time

This is the most basic method and is the first one that beginners will use.

- RGui (Microsoft ® Windows)
- Terminal
- Emacs/ESS, XEmacs
- JGR

1.2.2 Multiple lines at a time

For longer programs (called *scripts*) there is too much code to write all at once at the command prompt. Furthermore, for longer scripts it is convenient to be able to only modify a certain piece of the script and run it again in R. Programs called *script editors* are specially designed to aid the communication and code writing process. They have all sorts of helpful features including R syntax highlighting, automatic code completion, delimiter matching, and dynamic help on the R functions as they are being written. Even more, they often have all of the text editing features of programs like Microsoft®Word. Lastly, most script editors are fully customizable in the sense that the user can customize the appearance of the interface to choose what colors to display, when to display them, and how to display them.

R Editor (Windows): In Microsoft® Windows, R Gui has its own built-in script editor, called R Editor. From the console window, select **File ▸ New Script**. A script window opens, and the lines of code can be written in the window. When satisfied with the code, the user highlights all of the commands and presses **Ctrl+R**. The commands are automatically run at once in R and the output is shown. To save the script for later, click **File ▸ Save as...** in R Editor. The script can be reopened later with **File ▸ Open Script...** in RGui. Note that R Editor does not have the fancy syntax highlighting that the others do.

R WinEdt: This option is coordinated with WinEdt for L^AT_EX and has additional features such as code highlighting, remote sourcing, and a ton of other things. However, one first needs to download and install a shareware version of another program, WinEdt, which is only free for a while – pop-up windows will eventually appear that ask for a registration code. R WinEdt is nevertheless a very fine choice if you already own WinEdt or are planning to purchase it in the near future.

Tinn R / Sciviews K: This one is completely free and has all of the above mentioned options and more. It is simple enough to use that the user can virtually begin working with it immediately after installation. But Tinn R proper is only available for Microsoft® Windows operating systems. If you are on MacOS or Linux, a comparable alternative is Sci-Views - Komodo Edit.

Emacs/ESS: Emacs is an all purpose text editor. It can do absolutely anything with respect to modifying, searching, editing, and manipulating, text. And if Emacs can't do it, then you can write a program that extends Emacs to do it. Once such extension is called ESS, which stands for *E*-macs *S*-peaks *S*-tatistics. With ESS a person can speak to R, do all of the tricks that the other script editors offer, and much, much, more. Please see the following for installation details, documentation, reference cards, and a whole lot more: <http://ess.r-project.org>. *Fair warning:* if you want to try Emacs and if you grew up with Microsoft® Windows or Macintosh, then you are going to need to relearn everything you thought you knew about computers your whole life. (Or, since Emacs is completely customizable, you can reconfigure Emacs to behave the way you want.) I have personally experienced this transformation and I will never go back.

JGR (read “Jaguar”): This one has the bells and whistles of RGui plus it is based on Java, so it works on multiple operating systems. It has its own script editor like R Editor but with additional features such as syntax highlighting and code-completion. If you do not use Microsoft® Windows (or even if you do) you definitely want to check out this one.

Kate, Bluefish, etc There are literally dozens of other text editors available, many of them free, and each has its own (dis)advantages. I only have mentioned the ones with which I have had substantial personal experience and have enjoyed at some point. Play around, and let me know what you find.

1.2.3 Graphical User Interfaces (GUIs)

By the word “GUI” I mean an interface in which the user communicates with R by way of points-and-clicks in a menu of some sort. Again, there are many, many options and I only mention ones that I have used and enjoyed. Some of the other more popular script editors can be downloaded from the R-Project website at http://www.sciviews.org/_rgui/. On the left side of the screen (under **Projects**) there are several choices available.

R Commander provides a point-and-click interface to many basic statistical tasks. It is called the “Commander” because every time one makes a selection from the menus, the code corresponding to the task is listed in the output window. One can take this code, copy-and-paste it to a text file, then re-run it again at a later time without the

R Commander's assistance. It is well suited for the introductory level. Rcmdr also allows for user-contributed "Plugins" which are separate packages on CRAN that add extra functionality to the Rcmdr package. The plugins are typically named with the prefix RcmdrPlugin to make them easy to identify in the CRAN package list. One such plugin is the RcmdrPlugin.IPSUR package which accompanies this text.

Poor Man's GUI is an alternative to the Rcmdr which is based on GTK instead of Tcl/Tk. It has been a while since I used it but I remember liking it very much when I did. One thing that stood out was that the user could drag-and-drop data sets for plots. See here for more information: <http://wiener.math.csi.cuny.edu/pmg/>.

Rattle is a data mining toolkit which was designed to manage/analyze very large data sets, but it provides enough other general functionality to merit mention here. See [?] for more information.

Deducer is relatively new and shows promise from what I have seen, but I have not actually used it in the classroom yet.

1.3 Basic R Operations and Concepts

The R developers have written an introductory document entitled "An Introduction to R". There is a sample session included which shows what basic interaction with R looks like. I recommend that all new users of R read that document, but bear in mind that there are concepts mentioned which will be unfamiliar to the beginner.

Below are some of the most basic operations that can be done with R. Almost every book about R begins with a section like the one below; look around to see all sorts of things that can be done at this most basic level.

1.3.1 Arithmetic

```
2 + 3          # add
4 * 5 / 6      # multiply and divide
7^8           # 7 to the 8th power

[1] 5
[1] 3.333333
[1] 5764801
```

Notice the comment character #. Anything typed after a # symbol is ignored by R. We know that 20/6 is a repeating decimal, but the above example shows only 7 digits. We can change the number of digits displayed with options:

```
options(digits = 16)
10/3                # see more digits
sqrt(2)             # square root
exp(1)              # Euler's constant, e
pi
options(digits = 7) # back to default

[1] 3.33333333333333
[1] 1.414213562373095
[1] 2.718281828459045
[1] 3.141592653589793
```

Note that it is possible to set `digits` up to 22, but setting them over 16 is not recommended (the extra significant digits are not necessarily reliable). Above notice the `sqrt` function for square roots and the `exp` function for powers of e , Euler's number.

1.3.2 Assignment, Object names, and Data types

It is often convenient to assign numbers and values to variables (objects) to be used later. The proper way to assign values to a variable is with the `<-` operator (with a space on either side). The `=` symbol works too, but it is recommended by the R masters to reserve `=` for specifying arguments to functions (discussed later). In this book we will follow their advice and use `<-` for assignment. Once a variable is assigned, its value can be printed by simply entering the variable name by itself.

```
x <- 7*41/pi      # don't see the calculated value
x                # take a look

[1] 91.35494
```

When choosing a variable name you can use letters, numbers, dots “.”, or underscore “_” characters. You cannot use mathematical operators, and a leading dot may not be followed by a number. Examples of valid names are: `x`, `x1`, `y.value`, and `!y_hat`. (More precisely, the set of allowable characters in object names depends on one's particular system and locale; see *An Introduction to R* for more discussion on this.)

Objects can be of many *types*, *modes*, and *classes*. At this level, it is not necessary to investigate all of the intricacies of the respective types, but there are some with which you need to become familiar:

integer: the values 0, ± 1 , ± 2 , ...; these are represented exactly by R.

double: real numbers (rational and irrational); these numbers are not represented exactly (save integers or fractions with a denominator that is a power of 2, see [?]).

character: elements that are wrapped with pairs of "=" or ';

logical: includes TRUE, FALSE, and NA (which are reserved words); the NA stands for “not available”, *i.e.*, a missing value.

You can determine an object’s type with the `typeof` function. In addition to the above, there is the complex data type:

```
sqrt(-1)           # isn't defined
sqrt(-1+0i)        # is defined
sqrt(as.complex(-1)) # same thing
(0 + 1i)^2          # should be -1
typeof((0 + 1i)^2)
```

```
[1] NaN
[1] 0+1i
[1] 0+1i
[1] -1+0i
[1] "complex"
```

Note that you can just type `(1i)^2` to get the same answer. The NaN stands for “not a number”; it is represented internally as double.

1.3.3 Vectors

All of this time we have been manipulating vectors of length 1. Now let us move to vectors with multiple entries.

Entering data vectors

The long way: If you would like to enter the data 74, 31, 95, 61, 76, 34, 23, 54, 96 into R, you may create a data vector with the `c` function (which is short for *concatenate*).

```
x <- c(74, 31, 95, 61, 76, 34, 23, 54, 96)
x
```

```
[1] 74 31 95 61 76 34 23 54 96
```

The elements of a vector are usually coerced by R to the the most general type of any of the elements, so if you do `c(1, "2")` then the result will be `c("1", "2")`.

A shorter way: : The `scan` method is useful when the data are stored somewhere else. For instance, you may type `x <- scan()` at the command prompt and R will display 1: to indicate that it is waiting for the first data value. Type a value and press Enter, at which

point R will display 2:, and so forth. Note that entering an empty line stops the scan. This method is especially handy when you have a column of values, say, stored in a text file or spreadsheet. You may copy and paste them all at the 1: prompt, and R will store all of the values instantly in the vector `x`.

Repeated data; regular patterns: the `seq` function will generate all sorts of sequences of numbers. It has the arguments `from`, `to`, `by`, and `length.out` which can be set in concert with one another. We will do a couple of examples to show you how it works.

```
seq(from = 1, to = 5)
seq(from = 2, by = -0.1, length.out = 4)
```

```
[1] 1 2 3 4 5
[1] 2.0 1.9 1.8 1.7
```

Note that we can get the first line much quicker with the colon operator.

```
1:5
```

```
[1] 1 2 3 4 5
```

The vector `LETTERS` has the 26 letters of the English alphabet in uppercase and `letters` has all of them in lowercase.

Indexing data vectors

Sometimes we do not want the whole vector, but just a piece of it. We can access the intermediate parts with the `[]` operator. Observe (with `x` defined above)

```
x[1]
x[2:4]
x[c(1, 3, 4, 8)]
x[-c(1, 3, 4, 8)]
```

```
[1] 74
[1] 31 95 61
[1] 74 95 61 54
[1] 31 76 34 23 96
```

Notice that we used the minus sign to specify those elements that we do *not* want.

```
LETTERS[1:5]
letters[-(6:24)]
```

```
[1] "A" "B" "C" "D" "E"
[1] "a" "b" "c" "d" "e" "y" "z"
```

1.3.4 Functions and Expressions

A function takes arguments as input and returns an object as output. There are functions to do all sorts of things. We show some examples below.

```
x <- 1:5
sum(x)
length(x)
min(x)
mean(x)      # sample mean
sd(x)        # sample standard deviation

[1] 15
[1] 5
[1] 1
[1] 3
[1] 1.581139
```

It will not be long before the user starts to wonder how a particular function is doing its job, and since R is open-source, anybody is free to look under the hood of a function to see how things are calculated. For detailed instructions see the article “Accessing the Sources” by Uwe Ligges [?]. In short:

Type the name of the function without any parentheses or arguments. If you are lucky then the code for the entire function will be printed, right there looking at you. For instance, suppose that we would like to see how the `intersect` function works:

```
intersect

function (x, ...)
UseMethod("intersect")
<environment: namespace:prob>
```

If instead it shows `UseMethod(something)` then you will need to choose the *class* of the object to be inputted and next look at the *method* that will be *dispatched* to the object. For instance, typing `rev` says

```
rev

function (x)
UseMethod("rev")
<environment: namespace:base>
```

The output is telling us that there are multiple methods associated with the `rev` function. To see what these are, type

```
methods(rev)
```

```
[1] rev.default      rev.dendrogram* rev.zoo  
[4] rev.zooreg*
```

Non-visible functions are asterisked

Now we learn that there are two different `rev(x)` functions, only one of which being chosen at each call depending on what `x` is. There is one for `dendrogram` objects and a `default` method for everything else. Simply type the name to see what each method does. For example, the `default` method can be viewed with

```
rev.default
```

```
function (x)  
if (length(x)) x[length(x):1L] else x  
<environment: namespace:base>
```

Some functions are hidden by a *namespace* (see An Introduction to R [?]), and are not visible on the first try. For example, if we try to look at the code for `wilcox.test` (see Chapter ??) we get the following:

```
wilcox.test  
methods(wilcox.test)
```

```
function (x, ...)  
UseMethod("wilcox.test")  
<environment: namespace:stats>  
[1] wilcox.test.default* wilcox.test.formula*
```

Non-visible functions are asterisked

If we were to try `wilcox.test.default` we would get a “not found” error, because it is hidden behind the namespace for the package `stats` (shown in the last line when we tried `wilcox.test`). In cases like these we prefix the package name to the front of the function name with three colons; the command `stats::wilcox.test.default` will show the source code, omitted here for brevity.

If it shows `.Internal(something)` or `.Primitive(something)`, then it will be necessary to download the source code of R (which is *not* a binary version with an `.exe` extension) and search inside the code there. See Ligges [?] for more discussion on this. An example is `exp`:

exp

```
function (x) .Primitive("exp")
```

Be warned that most of the `.Internal` functions are written in other computer languages which the beginner may not understand, at least initially.

1.4 Getting Help

When you are using R, it will not take long before you find yourself needing help. Fortunately, R has extensive help resources and you should immediately become familiar with them. Begin by clicking **Help** on RGui. The following options are available.

Console: gives useful shortcuts, for instance, `Ctrl+L`, to clear the R console screen.

FAQ on R: frequently asked questions concerning general R operation.

FAQ on R for Windows: frequently asked questions about R, tailored to the Microsoft Windows operating system.

Manuals: technical manuals about all features of the R system including installation, the complete language definition, and add-on packages.

R functions (text)...: use this if you know the *exact* name of the function you want to know more about, for example, `mean` or `plot`. Typing `mean` in the window is equivalent to typing `help("mean")` at the command line, or more simply, `?mean`. Note that this method only works if the function of interest is contained in a package that is already loaded into the search path with `library`.

HTML Help: use this to browse the manuals with point-and-click links. It also has a Search Engine & Keywords for searching the help page titles, with point-and-click links for the search results. This is possibly the best help method for beginners. It can be started from the command line with the command `help.start()`.

Search help ...: use this if you do not know the exact name of the function of interest, or if the function is in a package that has not been loaded yet. For example, you may enter `plo` and a text window will return listing all the help files with an alias, concept, or title matching `'=plo='` using regular expression matching; it is equivalent to typing `help.search("plo")` at the command line. The advantage is that you do not need to know the exact name of the function; the disadvantage is that you cannot point-and-click the results. Therefore, one may wish to use the HTML Help search engine instead. An equivalent way is `??plo` at the command line.

search.r-project.org ...: this will search for words in help lists and email archives of the R Project. It can be very useful for finding other questions that other users have asked.

Apropos ...: use this for more sophisticated partial name matching of functions. See `?apropos` for details.

On the help pages for a function there are sometimes “Examples” listed at the bottom of the page, which will work if copy-pasted at the command line (unless marked otherwise). The `example` function will run the code automatically, skipping the intermediate step. For instance, we may try `example(mean)` to see a few examples of how the `mean` function works.

1.4.1 R Help Mailing Lists

There are several mailing lists associated with R, and there is a huge community of people that read and answer questions related to R. See [here](#) for an idea of what is available. Particularly pay attention to the bottom of the page which lists several special interest groups (SIGs) related to R.

Bear in mind that R is free software, which means that it was written by volunteers, and the people that frequent the mailing lists are also volunteers who are not paid by customer support fees. Consequently, if you want to use the mailing lists for free advice then you must adhere to some basic etiquette, or else you may not get a reply, or even worse, you may receive a reply which is a bit less cordial than you are used to. Below are a few considerations:

1. Read the [FAQ](#). Note that there are different FAQs for different operating systems. You should read these now, even without a question at the moment, to learn a lot about the idiosyncrasies of R.
2. Search the archives. Even if your question is not a FAQ, there is a very high likelihood that your question has been asked before on the mailing list. If you want to know about topic `foo`, then you can do `RSiteSearch("foo")` to search the mailing list archives (and the online help) for it.
3. Do a Google search and an `RSeek.org` search.

If your question is not a FAQ, has not been asked on R-help before, and does not yield to a Google (or alternative) search, then, and only then, should you even consider writing to R-help. Below are a few additional considerations.

- Read the [posting guide](#) before posting. This will save you a lot of trouble and pain.

- Get rid of the command prompts (>) from output. Readers of your message will take the text from your mail and copy-paste into an R session. If you make the readers' job easier then it will increase the likelihood of a response.
- Questions are often related to a specific data set, and the best way to communicate the data is with a `dump` command. For instance, if your question involves data stored in a vector `x`, you can type `dump("x", "")` at the command prompt and copy-paste the output into the body of your email message. Then the reader may easily copy-paste the message from your email into R and `x` will be available to him/her.
- Sometimes the answer the question is related to the operating system used, the attached packages, or the exact version of R being used. The `sessionInfo()` command collects all of this information to be copy-pasted into an email (and the Posting Guide requests this information). See Appendix ?? for an example.

1.5 External Resources

There is a mountain of information on the Internet about R. Below are a few of the important ones.

- The R- Project for Statistical Computing:: Go [there](#) first.
- The Comprehensive R Archive Network:: [That is where](#) R is stored along with thousands of contributed packages. There are also loads of contributed information (books, tutorials, *etc.*). There are mirrors all over the world with duplicate information.
- R-Forge:: [This is another location](#) where R packages are stored. Here you can find development code which has not yet been released to CRAN.
- R-Wiki:: There are many tips, tricks, and general advice [listed here](#). If you find a trick of your own, login and share it with the world.
- Other: the [R Graph Gallery](#) and [R Graphical Manual](#) have literally thousands of graphs to peruse. [R Seek](#) is a search engine based on Google specifically tailored for R queries.

1.6 Other Tips

It is unnecessary to retype commands repeatedly, since R remembers what you have recently entered on the command line. On the Microsoft® Windows R Gui, to cycle through the previous commands just push the ↑ (up arrow) key. On Emacs/ESS the command is

M-p (which means hold down the Alt button and press “p”). More generally, the command `history()` will show a whole list of recently entered commands.

- To find out what all variables are in the current work environment, use the commands `objects()` or `ls()`. These list all available objects in the workspace. If you wish to remove one or more variables, use `remove(var1, var2, var3)`, or more simply use `rm(var1, var2, var3)`, and to remove all objects use `rm(list = ls())`.
- Another use of `scan` is when you have a long list of numbers (separated by spaces or on different lines) already typed somewhere else, say in a text file. To enter all the data in one fell swoop, first highlight and copy the list of numbers to the Clipboard with Edit ▸ Copy (or by right-clicking and selecting Copy). Next type the `x <- scan()` command in the R console, and paste the numbers at the 1: prompt with Edit ▸ Paste. All of the numbers will automatically be entered into the vector `x`.
- The command `Ctrl+l` clears the display in the Microsoft® Windows R Gui. In Emacs/ESS, press `Ctrl+l` repeatedly to cycle point (the place where the cursor is) to the bottom, middle, and top of the display.
- Once you use R for awhile there may be some commands that you wish to run automatically whenever R starts. These commands may be saved in a file called `Rprofile.site` which is usually in the `etc` folder, which lives in the R home directory (which on Microsoft® Windows usually is `C:\Program Files\R`). Alternatively, you can make a file `.Rprofile` to be stored in the user’s home directory, or anywhere R is invoked. This allows for multiple configurations for different projects or users. See “Customizing the Environment” of *An Introduction to R* for more details.
- When exiting R the user is given the option to “save the workspace”. I recommend that beginners DO NOT save the workspace when quitting. If Yes is selected, then all of the objects and data currently in R’s memory is saved in a file located in the working directory called `.RData`. This file is then automatically loaded the next time R starts (in which case R will say `[previously saved workspace restored]`). This is a valuable feature for experienced users of R, but I find that it causes more trouble than it saves with beginners.

1.7 Exercises