

A R GUIDE

Version 0.07.12.30

Samuele Carcagno
sam.carcagno@gmail.com

Copyright ©2006–08 Samuele Carcagno.

Permission is granted to copy, distribute and/or modify this document under the terms of the GNU Free Documentation License, Version 1.2 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”.

Disclaimer: This document comes with NO WARRANTY whatsoever of being correct in any of its parts. This document is work in progress, check out http://samcarcagno.altervista.org/soft/carcagno_R_guide.pdf to download the latest version. This document will reach version number 1 when I think it's decent and complete enough. For the moment version number is 0, followed by year.month.day of last modification.

Contents

1	Installing R	1
1.1	Installing the Base Program	1
1.2	Installing Add-On Packages	2
2	A Simple Introduction to R	3
2.1	Firing up and Quitting R	3
2.2	Starting to Work with R	4
2.3	Getting Help	5
2.4	Working with a Graphical User Interface	7
3	Organising a Working Session	9
3.1	Setting and Changing the Working Directory	9
3.2	Objects	9
3.3	Saving and Using the “Workspace Image”	10
3.4	Working in Batch Mode	11
4	Data Types and Data Manipulation	12
4.1	Vectors	12
4.1.1	Indexing Vectors	13
4.1.2	The <code>seq()</code> Function	14
4.1.3	The <code>rep()</code> function	15
4.2	Matrix Facilities	15
4.2.1	Matrix Operations	17
4.3	Lists	18
4.4	Dataframes	19
4.5	Factors	22
4.6	Getting Info on R Objects	23
4.7	Changing the Format of Your Data	24
4.7.1	The “One Row per Observation” Format	25
4.7.2	The <code>stack</code> and <code>unstack</code> functions	29
4.7.3	<code>tapply</code> and <code>aggregate</code>	30
4.8	Creating and Editing Data Objects through a Visual Interface	31
4.9	Visualising Your Data	32

4.10	Reading Numbers in Exponential Notation	33
5	File Input/Output	34
5.1	Reading In Data from a File	34
5.1.1	<code>read.table()</code>	34
5.1.2	<code>scan()</code>	34
5.1.3	Low-Level File Input	36
5.2	Writing Data to a file	36
5.2.1	The <code>write.table</code> function	36
5.2.2	Low-Level File Output	37
6	Descriptive	38
6.1	Tables	38
6.2	The <code>scale()</code> Function	38
7	Graphics	39
7.1	Overview of R Graphics	39
7.2	The <code>plot()</code> Function	39
7.3	Drawing Functions	42
7.3.1	<code>matplot</code>	42
7.4	Barplots	42
7.5	Boxplots	47
7.6	Histograms	49
7.7	Stripcharts	49
7.8	Interaction Plots	51
7.9	Setting Graphics Parameters	54
7.9.1	Line type <code>lty</code>	55
7.9.2	Symbols with <code>pch</code>	57
7.10	Adding Elements to a Plot	58
7.10.1	Adding a Legend	58
7.10.2	Adding Text	59
7.10.3	Adding a Grid	59
7.10.4	Setting the Axes	60
7.11	Creating Layouts for Multiple Graphs	60
7.11.1	<code>mfrow</code> and <code>mfc</code> ol	60
7.11.2	<code>layout</code>	61
7.12	Graphics Device Regions and Coordinates	63
7.13	Plotting from Scratch	68
7.14	Colours for Graphics	69
7.15	Managing Graphic Devices	72
7.15.1	Opening Another Graphics Window	72
7.15.2	Exporting Graphics	72

8	Lattice or Trellis Graphics	73
8.1	Overview of Lattice Graphics	74
8.2	Introduction to model formulae and multi-panel conditioning	74
8.3	barchart	78
8.4	dotplot	80
8.5	histogram	80
8.6	Interaction Plots	80
8.7	Customising Lattice Graphics	81
8.7.1	Textual Elements	81
8.8	Writing Panel Functions	81
8.8.1	Combining Panel Functions	81
9	Probability Distribution Functions	85
9.1	The Bernoulli distribution	85
9.2	The binomial distribution	86
9.3	The normal distribution	86
10	Hypothesis Testing	90
10.1	χ^2 test	90
10.2	Student's t test	93
10.2.1	One sample t-test	93
10.2.2	Two samples t-test	93
10.3	The Levene Test for Homogeneity of Variances	94
11	Correlation and Regression	97
11.1	Linear Regression	98
12	ANOVA	99
12.1	One-Way ANOVA	99
12.1.1	One-Way ANOVA with Unequal Sample Sizes	101
12.2	Repeated Measures ANOVA	102
12.2.1	One Within Subjects Factor	102
12.2.2	Two Within Subjects Factors	105
12.2.3	Two Within Subjects Factors and One Between	108
13	How to Adjust the p-values for Multiple Comparisons	111
14	R Programming	113
14.1	Control Structures	113
14.2	String Processing	115
14.2.1	Using Regular Expressions	116
14.3	Tips and Tricks	117
14.3.1	Convert a String into a Command	117
14.4	Creating Simple R Packages	117
14.4.1	The Bare Minimum to Create a Package	118

15 Environment Customisation	122
16 ESS: Using Emacs Speaks Statistics with R	124
16.1 Using ESS for Editing and Debugging R Source Files	124
16.2 Using ESS to Interact with a R Process	126
17 Using Sweave to Write Documents with R and L^AT_EX	127
17.1 What's this Sweave?	127
17.2 Usage	128
18 Sound Processing	129
18.1 Libraries for Sound Analysis and Signal Processing	129
A Partial List of Packages by Category	130
A.1 Graphics Packages	130
A.2 GUI Packages	131
B Miscellaneous commands	132
C Other Manuals and Sources of Information on R	133
D Other Useful Statistics Resources	135
E Full Colours Table	136
F GNU Free Documentation License	143
1. APPLICABILITY AND DEFINITIONS	143
2. VERBATIM COPYING	145
3. COPYING IN QUANTITY	145
4. MODIFICATIONS	146
5. COMBINING DOCUMENTS	148
6. COLLECTIONS OF DOCUMENTS	149
7. AGGREGATION WITH INDEPENDENT WORKS	149
8. TRANSLATION	149
9. TERMINATION	150
10. FUTURE REVISIONS OF THIS LICENSE	150
ADDENDUM: How to use this License for your documents	150
References	152
Index	154

Chapter 1

Installing R

The information provided in this section is quite generic and limited, for detailed information, please refer to the “**R Installation and Administration Manual**” available at the CRAN website <http://cran.r-project.org/> CRAN stands for “Comprehensive R Archive Network”, and is the main point of reference for the R software, where you can find R sources, binaries, documentation, and add-on packages.

1.1 Installing the Base Program

First of all you need to install the base R program. This already implements most of the statistical functions you will ever need. There are two ways of getting it, you can either compile the source code yourself, or you can fetch the precompiled binaries for your specific operative system. The second way is the easiest one, and usually you will want to go with it.

GNU/Linux and Other Unices

Precompiled binaries are available for some GNU/Linux distributions, there is a list of these on the R FAQ. For other distributions you can build R from source. There are precompiled binaries for **Debian GNU/Linux**, so you can install the R base system and a number of add-on packages with the usual methods under Debian, that is, `apt-get`, `Synaptic` or whatever else you use.

Windows

There are precompiled binaries for Windows, you just have to download them, then double click on the installer’s icon to start the installation. Most Windows versions are supported.

Mac Os X

Precompiled binaries are also available for Mac Os X, you can download them and then double click on the installer's icon to start the installation.

1.2 Installing Add-On Packages

There is a growing number of packages that implement statistical functions that are not available with the base program (though you can pretty easily write your own functions with R), or provide other utilities. They're not strictly necessary, but if you keep using R, sooner or later you will end up wanting to install some of these packages.

GNU/Linux and Other Unices

There are different ways to get packages installed. For Debian there are precompiled versions of some packages, so you can get them with `apt-get` or whatever else you usually use to install Debian packages, this will also take care of possible dependencies (some packages need other packages to be installed in order to work). For other packages, a way to install them is to download the related tarball and then from a root console issue the command:

```
R CMD INSTALL /packagename
```

where "packagename" is the full path to the tarball you have downloaded. There are other ways and other specific options to install add-on packages. Again, refer to the "R Installation and Administration Manual" for further information.

Windows


The R's GUI on Windows provides an interface to download and install the packages directly from the internet. You can also install packages from a local repository, this is necessary for example if you don't have an internet connection from the computer you want to install the package on. In this case, you can download the package you want from another computer as a zip file, and then transfer it to your computer. At this point to install the package don't unzip it, just start the R GUI and click on etc

Chapter 2




A Simple Introduction to R

This chapter will give a simple introduction to R, just to get familiar with it and get a general idea of how it works. This chapter assumes no previous knowledge of programming or anything similar. If you know other computer languages, or have even a basic knowledge of programming, getting started will be easy. If you don't, don't worry, R syntax is very elegant and simple, it might take a little while, but after looking at some examples, and importantly, trying them out yourself, you'll be up and running without problems. This tutorial deals only with learning to use R from the command line, if you'd rather use R with a GUI, that is, with a "point and click" interface, please have a look at section 2.4 for some information on how to get started.

2.1 Firing up and Quitting R

Under GNU/Linux systems you can start R from a shell, just type `R` and press . Under Windows you can click on the R icon to start the R GUI.

You can save yourself a lot of typing by using the up arrow key  to retrieve past commands.

Commands can be terminated either by a semi-colon `;` or by pressing  and starting a newline. If you start a newline before a command is complete, R will prompt you to complete the command with a `+` sign, you can then complete the command. If you don't know how to complete the command and get stuck, you can stop R prompting you with the plus sign by pressing the  and  keys simultaneously.

To quit R type `quit()` or `q()`¹, R will ask you if you want to save the current session, if you answer `'y'`, R will save all the objects active in the current session and the command history.

¹On a Unix terminal you can also press `Ctrl-d` to exit

2.2 Starting to Work with R

The first thing you can try, is doing some math, at the command prompt type `5+4` and press `Enter`, the result will be

```
> 5+4
[1] 9
```

well, obviously 9. Other arithmetic operators are listed in Table 2.1.

Table 2.1: Arithmetic Operators

Symbol	Function
+	Addition
-	Subtraction
*	Multiplication
/	Division
^	Exponentiation

Now let's create a variable, we'll call it `foo`, and assign to it a number

```
> foo <- 5
```

the arrow '`<-`' is the assignment operator in R², in the above case it means the value of `foo` is 5, `foo` is an object, since it is of a numeric type, we can perform arithmetic operations on it:

```
> foo * 2
[1] 10
```

we can also store the results of an arithmetic operation in an object:

```
> another_foo <- 5^2
```

if you want to display the value of this new object you can use

```
> print(another_foo)
[1] 25
```

or, for short, just type its name and press `Enter`

```
> another_foo
[1] 25
```

²You can also use the symbol `=` as an assignment operator, however the use of the arrow symbol is more straightforward, in that it clearly shows the direction of the assignment and doesn't create confusion with the `==` (equal to) logical operator

since the two objects we have created are both numeric we can also do

```
> foo + another_foo
[1] 30
```

Let's look at something more interesting, we can create an object that stores a series of numbers, for example the money we have spent each day of a week, in Euros, we can do this using the `c()` function, which concatenates a series of values in a *vector*

```
> expenses <- c(7,8,15,20,9,45,3)
```

you might want to find out how much you've spent in average during the week, this is easily accomplished with the function `mean()`

```
> mean(expenses)
[1] 15.28571
```

the function `sd()` gives you the standard deviation

```
> sd(expenses)
[1] 14.24446
```

Surely you've wondered what the hell is that `[1]` that appears every time R gives you a result, now that we've introduced the vector we can get to it. Try to create a long vector, you can easily do this by creating a sequence of numbers, for example

```
> long_vec <- c(1:100)
```

will create a vector containing the sequence of numbers from 1 to 100, now try to display it and see what happens. All the elements of the vector won't fit in a single line of the screen, and at the start of each line you'll get between `[]` the index, that is the position, of the first element on that line. There's also a shorthand to create such a vector

```
> long_vec <- 1:100
```

2.3 Getting Help

R comes with an excellent online help facility which documents and gives examples for all available functions. There is also a web interface for the help system which is easier to use, you can start it with

```
> help.start()
```

this fires up a web browser from which you can access a search engine for all the available documentation. The documentation is also available as a pdf file, the "Full Reference Manual" which documents the base system. Printable pdf manuals are also available for all the other additional packages.

The Online Help System

You can quickly look up the documentation for a function, for example `sd`, with


```
> ?sd
```

or

```
> help(sd)
```

it is often indifferent using quotes or not, but sometimes they are required, for example

```
> ?* ## doesn't work
> Error: syntax error
> ?"*" ## this works!
```

to quit the help screen press  .

You can easily run the example code given in the help pages for a given function with

```
> example(function_name)
```

it is better to set the graphics parameter `ask` as `TRUE` before running the examples

```
> par(ask=TRUE)
```

to pause between successive plots, if there is more than one.

The online help system searches by default the available documentation for the base system and all the packages that are currently loaded. If you want to look the documentation for a function present in a package that is not loaded, you need to specify the package in question:

```
> help("levene.test", package=car)
```

if you know the function exists but don't know the package it is in, try

```
> help(levene.test, try.all.packages=TRUE)
Help for topic 'levene.test' is not in any loaded package but can be
found in the following packages:
```

Package	Library
car	/usr/lib/R/site-library

The function `help.search()` can be used when you don't exactly know the name of the function you're looking for

```
> help.search("levene")
Help files with alias or concept or title matching 'levene' using
fuzzy matching:

levene.test(car)          Levene's Test
...
```

2.4 Working with a Graphical User Interface

The version of R for Windows comes with a very limited graphical user interface (GUI), while the GNU/Linux version comes with no GUI at all. There are however several independent projects aimed at developing a GUI for R. The following sections give some information on them. An extensive GUI for R might be useful for beginners, however I'd recommend learning to use R from the command line, because it is a very elegant language and its syntax is pretty much easy to learn. Besides, working from the command line is faster.

The R Commander

An extensive GUI for R is provided by the `Rcmdr` package (R commander). This GUI allows you to do most of the operations you can do using R from the command line, through a point and click visual interface. The R commander is just a R package, so in order to use it, you need to have R installed in the first place, then you have to install the `Rcmdr` package and all the other packages it depends on. After everything is installed correctly, fire up R and call the R commander as you do with any other package:

```
>library(Rcmdr)
```

you will be greeted by a GUI with menus that allow you to type in data, perform statistical analyses and create graphs. The R commander works on both GNU/Linux and Windows platforms. For further information, please refer to the R commander manual or look up the following web page:

<http://socserv.socsci.mcmaster.ca/jfox/Misc/Rcmdr/>

Sciviews-R

Sciviews-R is a project aimed at providing an user friendly GUI for R. Currently it is only available for Windows. It provides different functionality from the R commander, but you can integrate the two together to get the best from each. Please, for further information refer to the Sciwies-R web page:

<http://www.sciviews.org/SciViews-R/index.html>.

JGR

The JGR package provides a clean, simple graphical user interface for R, which is platform independent. The package is written in JAVA and requires the JAVA SDK to run. More info at the developer's web-page:

<http://rosuda.org/JGR/index.shtml>

Chapter 3

Organising a Working Session

3.1 Setting and Changing the Working Directory

The command `getwd()` displays the pathname of the current working directory, that is where R will look for and store files if not otherwise instructed.

To change the current working directory, use the command `setwd("dirname")`, where `dirname` is the pathname of the working directory you're moving into. Note that this has to be an existent working directory, because R cannot create a new directory with this command. Here's an example of how to specify the pathname:

```
> setwd("C:/mydata/rats")
```

note that you have to use a slash "/" and not backslash "\" like you usually do in Windows to specify the pathname. You can also specify a pathname relative to your current working directory, without specifying the full pathname. It is indifferent using single ' ' or double " " quotes, this holds true when you need to quote character strings.

If you want to see the files present in your current working directory, use the command `dir()`.

It is also possible to issue commands to the OS from within R with the `system` function, for example

```
> system("ls")
```

under GNU/Linux or Unix systems, will list the files present in the current directory.

3.2 Objects

All the variables, functions, arrays etc. you work with in R are stored and manipulated as *objects*. To list all the objects currently active in your

workspace, you can use the command:

```
> objects()
```

or alternatively

```
> ls()
```

if you want to remove some of these objects from memory, you can use the command:

```
> rm(X, Y, W, foo, rats)
```

you can also give the variables to be removed, as a character vector:

```
rm(list=c("X","Y","foo","rats"))
```

if you want to remove all the objects in your workspace, you can combine the `ls()` and `rm()` commands as follows:

```
> rm(list=ls())
```

this however doesn't remove objects whose name starts with a dot, to remove also those you can use:

```
> rm(list=ls(all=TRUE))
```

3.3 Saving and Using the “Workspace Image”

You can use a “workspace image” that you have previously saved by starting R from the directory in which it was saved. In this way you can use the objects created in a previous session and the up arrow as well to retrieve commands from that session. To take full advantage from workspace images you'd better use different working directories for different analyses, studies, experiments and so on, in this way you can restore the workspace image of a specific analysis you were running and above all, you avoid accidentally overwriting objects from different analyses by creating another object with the same name during your current analysis.

When saving the workspace image R stores two files in the current working directory, one with the objects and one with the command history. This files begin with a dot under GNU/Linux and so are hidden.

You can save the workspace image either on exit, answering yes to the prompt you're given, or during a session with the command `save.image()`, the latter is a good measure against accidental losses of objects due to a power failure.

3.4 Working in Batch Mode

Executing commands written in a file from an R session

Instead of writing and executing commands line by line, it is often convenient to write the commands in a text file and then run them all at once in batch mode. You just write the commands with a text editor in a file, as if it were on the R console, save it in a directory, and then from within an R session issue the command:

```
> source("C:mydata/myfile.txt")
```

By default R displays only the results of the commands written in the source file, you can change this using the option:

```
> source("C:mydata/myfile.txt", echo=TRUE)
```

Executing commands written in a file from a shell

It is also possible to execute R commands written in a file without starting an R session: From within your system's shell (for example bash on GNU/Linux or dos on Windows) issue the command

```
$ R CMD BATCH myfile.r
```

where 'myfile.r' is the file you've written the R commands in.

Chapter 4

Data Types and Data Manipulation

4.1 Vectors

One of the simplest and among the most important data types in R is the vector, which can be numerical or containing strings of characters. A simple way to build a vector is through the function `c()`, which concatenates a series of data, for example:

```
> temperature <- c(34, 45, 23, 29, 26, 31, 44, 32, 19, 22, 34)
```

in this case the `c()` concatenates a series of numerical data into a vector and the assignment operator `<-` assigns it to the variable “temperature”, so that it can be retrieved later. Once the variable is created you can apply functions to it, for example

```
> mean(temperature)
[1] 30.81818
```

will compute the mean of the data vector. If you want to save the result of this function, you just have to assign it to another object

```
> mean_temp <- mean(temperature)
```

notice that in this case the value is assigned to the object `mean_temp` but it is not printed, you can display it with

```
> print(mean_temp)
[1] 30.81818
```

or for short, just calling the object

```
> mean_temp
[1] 30.81818
```

You can also perform simple arithmetic operations on a vector, for example:

```
> temperature + 10
[1] 44 55 33 39 36 41 54 42 29 32 44
```

will add 10 to each element of the vector.

You can also build vectors of characters, quoting each element of the vector

```
> colour <- c("blue", "green", "red" )
```

The `length()` function is used to access the number of element present in a vector

```
> length(colour)
[1] 3
```

4.1.1 Indexing Vectors

It is possible to access only subsets of data in a vector and also assign them to another vector. The most basic form of indexing is based on the position of the data in the vector. For example, to access only the datum in the third position of a vector called `temperature`, you would simply type:

```
> temperature[3]
```

if you would like to access the data in more than one position of the vector, let's say the first, the third and the sixth, you can again use the function concatenate inside the indexing command:

```
> temperature[c(1, 3, 6)]
```

to access the data from, say, the third position to the tenth position you can use:

```
> temperature[3:10]
```

and if you want to assign this subset to another vector called `"white"`, you can just type:

```
> white <- temperature[3:10]
```

if you want to access all the vector but the first five positions:

```
> temperature[-(1:5)]
```

since in R there is not a “delete” command, you can use this form of subsetting to remove elements of a vector, for example, if you would like to cancel the fourth element of the `temperature` vector you would write:

```
> temperature <- temperature[-4]
```

Furthermore, to access subsets of data you can do much more magic using logical operators and other tricks, for example if you want to access in a vector only the data greater than a certain value, you can use the > (greater than) logical operator:

```
> temperature[temperature>30]
```

In order to concatenate logical commands, you can use the & (and) logical operator:

```
> temperature[temperature>30 & temperature<35]
```

Table 4.1: Logical Operators

Operator	Description
&	Intersection (“and”)
&&	“and” (lazy evaluation)
	Union (“or”)
	“or” (lazy evaluation)
!	Negation
xor()	exclusive “or”
isTRUE(x)	

Table 4.2: Relational Operators

Operator	Description
<	Minor than
<=	Minor than or equal to
>	Greater than
>=	Greater than or equal to
==	Equal to
!=	Not equal to

It is also possible to apply labels to the positions of a vector, and then access the datum in a given position through its label:

```
> temperature <- c(34, 45, 23, 29, 26)
> names(temperature) <-c("Johnny", "Jack", "Tony", "Pippo", "Linda")
> temperature["Tony"]
Tony
  23
```

4.1.2 The seq() Function

The seq() function can be used to create evenly spaced sequences of numbers

```
> seq(1,10)
[1] 1 2 3 4 5 6 7 8 9 10
```

the default increment is 1, but you can change it with the option by

```
> seq(1, 1.9, by=0.1)
[1] 1.0 1.1 1.2 1.3 1.4 1.5 1.6 1.7 1.8 1.9
```

There's a shortcut for sequences with an increment of 1

```
> a <- 1:10
> a
[1] 1 2 3 4 5 6 7 8 9 10
```

4.1.3 The `rep()` function

You can use the `rep` function to create vectors which contain repetitions of the same elements. Let's start from the most simple use:

```
> vector1 <- rep(3, 13)
```

simply creates a vector of 13 elements, all having the value 3. More interestingly, you can repeat sequences of numbers:

```
> rep(1:4, 3)
[1] 1 2 3 4 1 2 3 4 1 2 3 4
```

as you see the above command repeats the sequence 1,2,3,4 three times. Furthermore, you can also specify the number of times a given element of the sequence should be repeated:

```
> rep(1:4, 3, each=2)
[1] 1 1 2 2 3 3 4 4 1 1 2 2 3 3 4 4 1 1 2 2 3 3 4 4
```

There are other ways to achieve this same effect, for example:

```
> rep(rep(1:4, c(2, 2, 2, 2)), 3)
[1] 1 1 2 2 3 3 4 4 1 1 2 2 3 3 4 4 1 1 2 2 3 3 4 4
```

would yield the same effect.

Even if it can look pretty useless at first, the `rep()` function comes in very handy for example, when you want to transform the data in a table from “one row per subject”, to “one row per observation”, which is necessary for example to run a repeated measures ANOVA with the `aov()` function. `rep()` makes it all easier as you can create vectors in which the occurrence of the levels of a factor are repeated over and over.

4.2 Matrix Facilities

There are different ways for creating a matrix in R, you often start from a vector, and then transform it into a matrix with the `matrix` function:

```

> matr <- c(3, 5, 6, 2, 5, 7, 9, 1, 5, 4, 2, 3)
> matr <- matrix(matr, ncol=3, byrow=TRUE)
> matr
      [,1] [,2] [,3]
[1,]    3    5    6
[2,]    2    5    7
[3,]    9    1    5
[4,]    4    2    3

```

you give to the `matrix` function either the `ncol` or the `nrow` parameters to specify the layout of the matrix. The default method that R uses to fill in the matrix is by columns, so if you want to fill it by rows, you need to set true the option `byrow`, as in the example above.

Matrix indexing is similar to vector indexing:

```

> matr[2,1]
[1] 2

```

the first index refers to the row number, and the second index to the column number. Omitting one of the two indexes is useful for slicing, for example

```

> matr[,2]
[1] 5 5 1 2

```

gives *all the rows* in the second column. This could alternatively been written as

```

> matr[1:4,2]

```

where the index is a *range* of rows. This notation is useful when you want to extract only part of the rows or columns, for example

```

> matr[2:4,2]
[1] 5 1 2

```

When the rows or columns to be extracted are not consecutive, you can use a *vector* of indexes for slicing

```

> matr[c(2,4),2]
[1] 5 2

```

To query the dimension of the matrix you can use the `dim` function, in this case our matrix has 4 rows and 3 columns:

```

> dim(matr)
[1] 4 3

```

The rows and columns of a matrix can also be assigned a name attribute through the `dimnames` function. The `dimnames` have to be a list of character vectors the same length as the matrix dimensions they refer to:

```

> dimnames(matr) <- list(c("row1", "row2", "row3", "row4"),
                        + c("col1", "col2", "col3"))
> dimnames(matr) #dimnames is a list
[[1]]
[1] "row1" "row2" "row3" "row4"

[[2]]
[1] "col1" "col2" "col3"
> matr #now the matrix is printed with its dimnames
      col1 col2 col3
row1    3    5    6
row2    2    5    7
row3    9    1    5
row4    4    2    3

```

the names of the rows can be retrieved with `rownames` and the names of the columns with `colnames`:

```

> rownames(matr)
[1] "row1" "row2" "row3" "row4"
> colnames(matr)
[1] "col1" "col2" "col3"

```

and these names can be indirectly used for sub-setting;

```

> matr[1, which(colnames(matr) == 'col2')]
[1] 5

```

4.2.1 Matrix Operations

The function `t()` gives the transpose of a matrix. The inverse of a matrix is obtained through the function `solve()`. Some other operators are listed in Table 4.3. Example:

```

> beta <- solve(t(x)%*%x)%*%(t(x)%*%y)

```

Table 4.3: Matrix Operations

Operator	Function
<code>%*%</code>	Matrix Multiplication
<code>det()</code>	Determinant
<code>solve()</code>	Inverse

4.3 Lists

Lists are objects that can contain elements of different modes (e.g numeric, character, logical), as well as other objects (vectors, matrices and also other lists). Let's build a small list to see how we can work on it:

```
> vec1 <- 1:12
> vec2 <- c('w','h','m')
> mylist <- list(vec1,vec2)
> mylist
[[1]]
 [1]  1  2  3  4  5  6  7  8  9 10 11 12

[[2]]
 [1] "w" "h" "m"
```

The syntax for subsetting a list is a bit awkward (but as we'll see later, naming the elements of a list makes things easier). To access an element of a list you can use the double brackets notation, for example

```
> mylist[[1]]
 [1]  1  2  3  4  5  6  7  8  9 10 11 12
```

returns the first element of the list 'mylist', which is a vector of length 12, if you want to access, say, the third element of this vector, the syntax is as follows

```
> mylist[[1]][3]
 [1] 3
```

Naming the elements of the list makes things easier

```
> mylist<-list(a=vec1,b=vec2)
> mylist
$a
 [1]  1  2  3  4  5  6  7  8  9 10 11 12

$b
 [1] "w" "h" "m"
```

now the first element of the list is named 'a', and the second 'b', and we can access them with a special "dollar sign" notation

```
> mylist$a ## return element of the list named a
 [1]  1  2  3  4  5  6  7  8  9 10 11 12
> mylist$a[1:3]
 [1] 1 2 3
```


It is also possible to use the double brackets notation with names

```
> myList[["a"]][3]
[1] 3
```

To eliminate an element of a list, set it to `NULL`

```
> vec1 <- c(1, 2, 3)
> vec2 <- c("a", "b", "c")
> myList <- list(vec1=vec1, vec2=vec2)
> myList
$vec1
[1] 1 2 3
```

```
$vec2
[1] "a" "b" "c"
```

```
> myList$vec1 <- NULL
> myList
$vec2
[1] "a" "b" "c"
```

notice that this is different from eliminating an element of the vectors contained in the list, you can do the latter with

```
> myList$vec2 <- myList$vec2[-1]
> myList
$vec2
[1] "b" "c"
```

4.4 Dataframes

Dataframes are one of the most important objects in R. You can think of it as a rectangular data structure, in which each column stores either the values of a numeric variable, or the levels of a factor, and each row represents an observation. Let's look at an example, we'll build a dataframe from 3 vectors, the first vector stores a variable, number of beers drunk during a week for twelve young people, the second is a factor vector, that specifies for each person whether he/she is a university student or not, so it has two levels, the third is also a factor vector, which tells the sex of each person, so it has two levels as well.

```
> n_beers <- c(6, 8, 4, 8, 9, 4, 5, 3, 4, 2, 3, 1)
> occupation <- rep(c("s", "w"), 6)
> sex <- c(rep("m", 6), rep("f", 6))
> occupation <- as.factor(occupation)
> sex <- as.factor(sex)
```

well, now let's create the dataframe

```
> dats <- data.frame(n_beers,occupation,sex)
```

it's as simple as this, you have just put the three vectors together, let's have a look at it

```
> dats
  n_beers occupation sex
1        6          s  m
2        8          w  m
3        4          s  m
4        8          w  m
5        9          s  m
6        4          w  m
7        5          s  f
8        3          w  f
9        4          s  f
10       2          w  f
11       3          s  f
12       1          w  f
```

as we said, each row holds the data of a single observation, in this case it corresponds to the data of a subject, but as we'll see later, this is not always necessarily true. Each row gives a full specification for each observation, we know that the first subject drunk 6 beers, he's a student, and he's male, and we could tell the same data for the other subjects. Since we have all this information, we could now compare for example the number of beers drunk by male vs females, or by male students vs female students. There are special functions to compute these values quickly like `tapply` and `by` (see sec. 4.7.3), and other functions to get statistical tests, they will be dealt with as we go along.

Accessing Parts of a Dataframe

You can access, or refer to a column of a dataframe with the `$` operator, in the example above suppose we removed all the original variables after creating the dataframe

```
> rm(n_beers,occupation,sex)
```

we can't now access them directly by name

```
> mean(n_beers)
Error in mean(n_beers) : Object "n_beers" not found
```

we have to retrieve them from the dataframe

```
> mean(dats$n_beers)
[1] 4.75
```

the example might seem artificial (why did I remove them in the first place?), but very often you read in the data directly as a dataframe with the `read.table` function (see sec. 5.1.1, so you'll have to access them from the dataframe. Another option is to use the function `attach`, which attaches the dataframe to the path that R searches when evaluating a variable, in this way you don't have to refer to the dataframe to access the values of a variable

```
> attach(dats)
> mean(n_beers)
[1] 4.75
```

this is OK only if you're working with a single dataframe, and you don't want to manipulate the variables in it. In fact if you accidentally attach two dataframes that share some variable names, or you try to change an object of a dataframe after it has been attached, strange things may happen, you've been warned, the details are in the R manual. To function `detach` detaches the dataframe from the search path.

Change the Names of Variables in a Dataframe

Sometimes you might want to change the names of the variables in a dataframe, for example when you create new dataframes with the `unstack` function, or just because you don't like the way you called it initially. You can visualise the names for the variables with the function `names`

```
> names(dats)
[1] "n_beers"      "occupation" "sex"
```

or if you want to see just the first one

```
> names(dats)[1]
[1] "n_beers"
```

you can change it with a simple assignment

```
> names(dats)[1] <- "beers"
```

or if you want to change more than one

```
> names(dats) <- c("brs", "occ", "sx")
```

Other Ways to Subset a Dataframe

A dataframe actually is just a special kind of list (a list of class dataframe), so we can use the normal list notation to subset dataframes

```
> dats[[1]]
[1] 6 8 4 8 9 4 5 3 4 2 3 1
```

Sometimes it's useful to think of a dataframe as a matrix, and use matrix notation for subsetting

```
> dats[1,] ## extract first row, all columns
  n_beers occupation sex
1       6          s   m
> dats[which(dats$n_beers>4),] ## records for subs who drink > 4 beers
  n_beers occupation sex
1       6          s   m
2       8          w   m
4       8          w   m
5       9          s   m
7       5          s   f
```

4.5 Factors

Data Vectors can be made not only of numerical values or of strings, but also *factors*. Factor vectors are very similar to character vectors, and could be seen as character vectors with some special properties. A factor vector usually consists of two or more levels, and can be created with the `factor` command. For example, suppose we are studying the drinking habits of 6 individuals, if one of the factors we are interested in is their sex, we can create a factor vector to encode it, if the first 3 are males and the other 3 females:

```
> sex <- factor(c("m", "m", "m", "f", "f", "f"))
> sex
[1] m m m f f f
Levels: f m
```

as you can see a factor has a “levels” attribute that specifies the possible values the factor can assume, and by default it is given by the unique values the factor vector can assume, sorted in alphabetical order.

In some cases, you might want to assign a label to each factor level. This is useful when you want to keep the level name short and simple, for ease of manipulation, but at the same time you would like to have it print nice labels in plots. When assigning a label it is best to also explicitly assign the factor levels in the same order, otherwise you might accidentally assign the wrong label to a given level:

```

> sex <- factor(c("m", "m", "m", "f", "f", "f"),
+ levels=c("m", "f"), labels=c("male", "female"))
> sex
[1] male   male   male   female female female
Levels: male female

```

In some cases it is also useful to set the factor as `ordered`, this might allow for example to have the levels of a factor plotted in a given order different from the alphabetical order

```

> sex <- factor(c("m", "m", "m", "f", "f", "f"),
+ levels=c("m", "f"), labels=c("male", "female"),
+ ordered=TRUE)

```

Creating Factors with `gl`

A handy function for creating factors for data with a regular pattern of factor levels is `gl`

```

> sex <- gl(2, 3, label=c("male", "female"), length=6)
> sex
[1] male   male   male   female female female

```

the first argument to the function specifies the number of levels, and the second argument the number of consecutive repetitions of each level, the pattern is repeated up to the number of elements specified by the `length` argument. Notice the different pattern created when the number of consecutive repetitions is set to 1 and the total length is left unchanged

```

> sex <- gl(2, 1, label=c("male", "female"), length=6)
> sex
[1] male   female male   female male   female
Levels: male female

```

4.6 Getting Info on R Objects

The most useful function to summarise information about many R object is `str()`:

```

> a <- seq(0, 10, .1)
> str(a)
 num [1:101] 0.0 0.1 0.2 0.3 0.4 ...
> b <- c("a", "b", "c")
> str(b)
 chr [1:3] "a" "b" "c"

```

in the example shown `str` gives information on the storage mode of the two vectors, numeric for `a` and character for `b`, it also gives information on the number of elements present and on the shape of the array (compare the output of `str` for a vector and a matrix). `str` gives also useful compact summaries of the contents of lists, including nested lists.

Another useful function is `mode`, which gives the storage mode of an R object:

```
> mode(a)
[1] "numeric"
> mode(b)
[1] "character"
```

4.7 Changing the Format of Your Data

In general statistical software require your data to be entered in a specific format in order to perform statistical analyses on them, and R is no exception. R provides many powerful functions to change the format of your data if they happen to be in a format that is not proper for applying a given statistical function on them. The process of changing the format of your data with these functions might seem very complicated at first, however you should keep in mind the following things:

- You don't really need to learn all of the functions that R provides to manipulate your data and change their format. Once you learn a procedure that does the job you can stick on it and you'll be fine most of the times.
- Once you understand the "logic" and the structure of the data format that R expects to apply some statistical functions, changing the layout of your data to match this structure will be easy. Moreover the data format R wants is most of the times one and only one: The "one row per observation format", which will be explained below.
- In this tutorial you will see different examples in which the format of the data, stored in a given file don't match the format R wants. This is just for illustrative purposes. In real life if you're doing a research or an experiment, you can often gather your data in a format that is already suitable for performing statistical analyses.
- You don't really have to use R to change the layout of your data if you don't like the functions it provides to do this job. You can always use some external programs to achieve the same results, for example spreadsheets programs such as **gnumeric** or the one provided with the **open office** suite. What's really important is that you understand the structure of the format the R expects.

This said, I would advice you to learn some of the functions that R provides to rearrange your data for at least two reasons: 1) They're very powerful and can actually save you time once you learn them, and 2) many examples in this tutorial and in other books use them, so you'll often need to know them to understand what's going on :)

4.7.1 The “One Row per Observation” Format

While statistic textbooks and scientific articles often show data in a format that is suitable and immediate for the “human eye”, like the one shown in table 4.4, statistical software often don't quite like it and would rather have the same data rearranged as shown in Table 4.5.

Table 4.4: Data Format Suitable for the “Human Eye”.

Group A	Group B	Group C
5	4	7
2	4	5
3	5	7
4	5	8
6	4	7

Table 4.5: Data Format Suitable for R.

Value	Group
5	A
2	A
3	A
4	A
6	A
4	B
4	B
5	B
5	B
4	B
7	C
5	C
7	C
8	C
7	C

The main difference is that while in the first format you have more than one observation in the same row, and you can identify the group to

which each observation belongs to through the column headers (“Group A”, “Group B” and “Group C”), in the second format you have only one observation for each row. This observation is then fully identified with a “label” that appears in the second column. A better way to describe the second column is to say that in the above case “Group” is a **factor**, and “A”, “B” and “C” define the **levels** of this factor for each group.

You could also be running an experiment in which you manipulate more than one factor. For example you might have two groups, “Patients” and “Controls”, which are tested under two conditions “1” and condition “2”. In this case you might display your data as shown in Table 4.6 to make them easily readable to humans. However for analysing your data with R, in this case you would need to add another column specifying the levels of the other factor for each observation as shown in Table 4.7

Table 4.6: Data Format Suitable for the “Human Eye” with More than One Factor.

Patients		Controls	
Condition 1	Condition 2	Condition 1	Condition 2
7	6	6	4
5	4	5	2
8	7	7	4
8	8	6	5
6	5	5	3

Finally, you might be running an experiment with a repeated measures design, in which all subjects are exposed to all the levels of all the within subjects factors. For example you might have your subjects recall word lists either under the effects of a drug or not (factor 1) and with words concrete or abstract words (factor 2). In this case the data for presentation might look like the ones in Table 4.8, in which each row represents a single subject. Again for R you need to rearrange the data so that each row represents a single observation, and in the case of a repeated measures design you need to add another column that identifies the levels of the “subjects” factor as shown in Table 4.9

Table 4.7: Data Format Suitable for R with More than One Factor.

Value	Group	Condition
7	P	1
5	P	1
8	P	1
8	P	1
6	P	1
6	P	2
4	P	2
7	P	2
8	P	2
5	P	2
6	C	1
5	C	1
7	C	1
6	C	1
5	C	1
4	C	2
2	C	2
4	C	2
5	C	2
3	C	2

Table 4.8: Data Format Suitable for the “Human Eye”, repeated measures design.

Drug		No-Drug	
Concrete	Abstract	Concrete	Abstract
7	6	6	4
5	4	5	2
8	7	7	4
8	8	6	5
6	5	5	3

Table 4.9: Data Format Suitable for R, repeated measures design.

Value	Drug Exposure	Word Type	Subject
7	D	C	1
5	D	C	2
8	D	C	3
8	D	C	4
6	D	C	5
6	D	A	1
4	D	A	2
7	D	A	3
8	D	A	4
5	D	A	5
6	N	C	1
5	N	C	2
7	N	C	3
6	N	C	4
5	N	C	5
4	N	A	1
2	N	A	2
4	N	A	3
5	N	A	4
3	N	A	5

4.7.2 The `stack` and `unstack` functions

One of the utilities provided by R to manipulate the format of your data is the `stack` function. If you have your data in a **dataframe** with a layout similar to that shown in Table 4.4, you can use the `stack` function to get a “one row per observation” format. What the `stack` function does is to create a single long vector from the vectors you have in your dataframe, and an additional factor vector which identifies the level for each observation. Here’s an example, the data are in the file `stack.txt`:

```
> dats<-read.table("stack.txt",header=TRUE)
```

this creates the dataframe

```
> dats <- stack(dats)
```

this reshapes the dataframe into a “one row per observation form”

Please note that R assigns names to the vectors in the new dataframe, you can see them in the header of the dataframe, you might need to know them for successive operations.

The `unstack` function simply does the opposite of the `stack` function, and you can use it if you want to switch back to your original dataframe format.

```
> dats <- unstack(dats)
```

The `unstack` function can do more tricks, if you have a dataframe with one observation per row, a column with a response variable, and two or more factor columns, you can unstack the values of the response variable according to the levels of only one factors or according to the levels of two or more factors. Suppose `lat` is your response variable, and you have two factors, `congr` with 3 levels and `isi`, also with 3 levels. The command:

```
> dats <- unstack(dats, form=lat~congr)
```

unstack the values of the response variable according to the levels of the `congr` factor, thus creating 3 columns, one for each level.

The command:

```
> dats <- unstack(dats,form=lat~congr:isi)
```

unstacks the values of the response variable according to the levels of both factors, thus creating $3 \times 3 = 9$ columns, the first column contains the values at level 1 of `congr` and level 1 of `isi`, the second one contains the values at level 1 of `congr` and level 2 of `isi`, and so on for all the possible combinations.

4.7.3 tapply and aggregate

The `tapply` function allows you to extract information from a dataframe, for example the mean or standard deviation of a given variable on the bases of one or more factor. The function name is related to the fact that it is used to *apply* a function (e.g. the mean) to a subsets of the dataframe chosen on the basis of one or more factors. We'll use the `InsectSprays` dataset to illustrate the use of `tapply`. The dataset contains the number of insects still alive in agricultural experimental units treated with six different types of pesticide.

```
> data(InsectSprays)
> head(InsectSprays)
  count spray
1     10    A
2      7    A
3     20    A
4     14    A
5     14    A
6     12    A
> meanSpray <- tapply(X=InsectSprays$count,
+                     INDEX=InsectSprays$spray,
+                     FUN=mean)
> meanSpray
      A      B      C      D      E      F
14.500000 15.333333  2.083333  4.916667  3.500000 16.666667
```

the arguments to `tapply` are `X`, the column of the dataframe to which the function should be applied, `INDEX`, the factor used for subsetting the dataframe, and `FUN`, the function to be applied. The function returns an array, in this case a vector, but can be a matrix, or multi-dimensional array. The `INDEX` argument indeed can be a *list* of factors, in this case the function chosen is applied to group of values given by a unique combination of the levels of these factors. We'll see an example by modifying the `InsectSprays` dataset including another fictitious factor. The new factor will be the season in which the fields were sprayed. The values will be returned in a matrix.

```
> season <- gl(4, 3, 72, labels=c("winter", "spring",
+                               "summer", "autumn"))
> Ins <- data.frame(InsectSprays, season)
> meanSpray <- tapply(X=Ins$count,
+                     INDEX=list(Ins$spray, Ins$season),
+                     FUN=mean)
> meanSpray
      winter  spring  summer  autumn
```

```

A 12.333333 13.333333 16.666667 15.666667
B 16.333333 13.666667 17.666667 13.666667
C  2.666667  2.000000  2.000000  1.666667
D  6.666667  4.333333  5.000000  3.666667
E  3.666667  4.666667  1.666667  4.000000
F 11.666667 17.666667 16.333333 21.000000

```

The `tapply` is often very useful, for example, after having calculated the means in this way, it is very easy to visualise the data with a barplot

```
> barplot(meanSpray, beside=TRUE, legend=T)
```

The `aggregate` function is very similar to `tapply`, but rather than returning an array, it returns a dataframe, which can be useful in some situations.

```

> InsDf <- aggregate(x=Ins$count,
+                   by=list(sprayType=Ins$spray, season=Ins$season),
+                   FUN=mean)
> InsDf
  sprayType season      x
1      A winter 12.333333
2      B winter 16.333333
3      C winter  2.666667
4      D winter  6.666667
5      E winter  3.666667
6      F winter 11.666667
7      A spring 13.333333
8      B spring 13.666667
.....

```

if you don't give a name to the grouping factors as we did with `SprayType=Ins$spray` a default name will be given. One slightly annoying thing is that, as far as I know, it is not possible to assign a name to the resulting variable (it will just be named `x`). However it can be changed afterwards, here's a solution that should work whatever the number of factors in the dataframe:

```

> names(InsDf)
[1] "sprayType" "season"    "x"
> names(InsDf)[which(names(InsDf)=="x")] <- "count"
> names(InsDf)
[1] "sprayType" "season"    "count"

```

4.8 Creating and Editing Data Objects through a Visual Interface

If you want to use a visual interface for creating a dataframe, first create an empty dataframe with:

```
> mydataframe <- data.frame()
```

then you can call a spreadsheet like editor to fill in the dataframe with:

```
> data.entry(mydataframe)
```

or

```
> fix(mydataframe)
```

If the data object is a vector, `fix` will call a text editor to edit the object instead of the spreadsheet like interface, so if you want the latter, use the function `data.entry` instead. However, using a simple text editor for fixing a vector might be more practical, if you want to use a different text editor from the one that `fix` calls by default, you can change the `editor` option:

```
> fix(myvector, editor="emacs")
```

or just call the editor on the object:

```
> emacs(myvector)
```

On Windows you could try:

```
> fix(myvector, editor="notepad")
```

4.9 Visualising Your Data

As you probably have already noticed, after you've created a data object, just typing its name on the console and pressing **Enter** will display the values it contains. Sometimes however, you might want to see only part of the data, for example to do some checking, or because the data object is too big and it's not printed nicely on the console. The functions `head` and `tail` let you look only at the first or the last part of your data respectively. For example, suppose `dats` is a dataframe, the command:

```
> head(dats)
```

will print only the first 6 observations. You can visualise more (or less) than 6 observations by setting the `n` option:

```
> head(dats, n=10)
```

4.10 Reading Numbers in Exponential Notation

Often R prints out numbers in exponential notation, this may lead to confusion. In order to understand the exponential notation is first necessary to introduce the *scientific notation*.

A number in scientific notation is in the form $a \cdot 10^b$, for example 300 could be written in scientific notation as $3 \cdot 10^2$. The components of a number in scientific notation are also named as *mantissa* $\cdot 10^{\text{characteristic}}$. Remember that a number with a negative exponent, for example 10^{-2} can be rewritten as

$$10^{-2} = \frac{1}{10^2} = 0.01$$

so, for example, 0.003 can be rewritten in scientific notation as $3 \cdot 10^{-3}$, because

$$3 \cdot 10^{-3} = 3 \cdot \frac{1}{10^3} = 0.003$$

R, as most calculators doesn't actually use the scientific notation, it uses instead the *exponential* notation. The exponential notation is a shorthand version of the scientific notation, in which, for example, 10^3 is replaced by $e3$, where e stands for *exponent*. So in our previous examples 300 would be written as $3e2$ and 0.003 would be written as $3e3$. Below are a few conversion examples.

Number	Scientific Notation	Exponential Notation
10	$1 \cdot 10^1$	$1e1$
20	$2 \cdot 10^1$	$2e1$
200	$2 \cdot 10^2$	$2e2$
350	$3.5 \cdot 10^2$	$3.5e2$
0.00353	$3.53 \cdot 10^{-3}$	$3.53e - 3$

As a quick and dirty rule, remember that when you're multiplying a number by 10^{exponent} , as you add to the exponent, you're adding a 0 to the number, or shifting the point by one position towards the right if it's a decimal number. As you subtract to the exponent, you're deleting a 0 to the number, or shifting the point by one position towards the left.

Chapter 5

File Input/Output

5.1 Reading In Data from a File

5.1.1 `read.table()`

If the data are in a table-like format, with each row corresponding to an observation or a single case, and each column to a variable, the most convenient function to load them in R is `read.table`. This function reads in the data file as a dataframe. For example the data in the data file “`ratsData.txt`” that contain information (height, weight, and species) of 6 rats can be easily read in with the following command:

```
> ratsData <- read.table(file="ratsData.txt", header=TRUE)
> ratsData
  identifier height weight species
1      sa01    3.2   300        A
2      sa02    2.6   246        A
3      sa03    2.9   317        A
4      sb01    2.4   229        B
5      sb02    2.5   230        B
6      sb03    2.4   245        B
```

in this case we’ve set the option `header=TRUE` because the file contains a header on the first line with the variable names.

5.1.2 `scan()`

Another very handy function for reading in data is `scan`, it can easily read in both tabular data in which all the columns are of the same type, or they are of different type, as long as they follow a regular pattern. We’ll read in the “`ratsData.txt`” file as an example:

```
> x <- scan(file="ratsData.txt", what=list(identifier=character()),
```



```

+           height=numeric(), weight=numeric(), species=character()),
+           skip=1)
Read 6 records
> str(x)
List of 4
 $ identifier: chr [1:6] "sa01" "sa02" "sa03" "sb01" ...
 $ height    : num [1:6] 3.2 2.6 2.9 2.4 2.5 2.4
 $ weight    : num [1:6] 300 246 317 229 230 245
 $ species   : chr [1:6] "A" "A" "A" "B" ...

```

the `what` argument tells the `scan` function the mode (numeric, character, etc...) of the elements to be read in, if `what` is a list of modes, then each corresponds to a column in the data file. So in the above example we have the first column, the variable `identifier`, which is of character mode, the second column, `height` is of numeric mode, like the third column, `weight`, while the last column, `species` is of character mode again. Notice that we're telling `scan` to skip the first line of the file (`skip=1`) because it contains the header. The object returned by `scan` in this case is a list, we can get the single elements of the list, corresponding to each column of the data file, with the usual methods for lists:

```

> x[[1]] #first item in list (first column in file)
[1] "sa01" "sa02" "sa03" "sb01" "sb02" "sb03"
> x[["weight"]] #item named weight in list (third column in file)
[1] "300" "246" "317" "229" "230" "245"

```

`scan` can do much more than what was shown in this example, like specifying the separator between the fields of the data file (comma, tabs, or whatever else, defaults to blank space), or specifying the maximum number of lines to be read, see `?scan` for more details. I'll present just another simple example in which we'll read in a file whose data are all numeric. The file is `rts.txt`

```

> x <- scan("rts.txt", what=numeric())
Read 36 items
> str(x)
 num [1:36] 0.12 0.132 0.102 0.096 0.103 0.087 0.113 ...

```

in this case the object returned is a long vector of the same mode as the `what` argument, a numeric vector. The file is however organised into three columns, which represent three different numeric variables. It is easy to reorganise our vector to reflect the structure of our data:

```

> nRows <- length(x)/3
> xm <- matrix(x, nrow=nRows, byrow=TRUE)
> xm
     [,1] [,2] [,3]

```

```

[1,] 0.120 0.132 0.102
[2,] 0.096 0.103 0.087
[3,] 0.113 0.134 0.109
[4,] 0.132 0.147 0.123
[5,] 0.124 0.139 0.124
[6,] 0.105 0.115 0.102
[7,] 0.109 0.129 0.097
[8,] 0.143 0.150 0.119
[9,] 0.127 0.145 0.113
[10,] 0.098 0.117 0.092
[11,] 0.115 0.126 0.098
[12,] 0.117 0.132 0.103

```

so we've got a matrix with the 3 columns of data originally found in the file, turning it into a dataframe would be equally easy at this point

```
xd <- as.data.frame(xm)
```

However also in this case it is possible to simply read in each column of the file as the element of a list:

```

> x <- scan("rts.txt", what=list(v1=numeric(),
+                               v2=numeric(), v3=numeric()))
Read 12 records
> str(x)
List of 3
 $ v1: num [1:12] 0.12 0.096 0.113 0.132 0.124 0.105 0.109 ...
 $ v2: num [1:12] 0.132 0.103 0.134 0.147 0.139 0.115 0.129 ...
 $ v3: num [1:12] 0.102 0.087 0.109 0.123 0.124 0.102 0.097 ...

```

5.1.3 Low-Level File Input

Sometimes the file to be read is not nicely organised into separate columns each representing a variable, in this case the function `readLines` can either read the file one line at the time, or all the lines at once, to be then further processed to extract the data.

5.2 Writing Data to a file

5.2.1 The `write.table` function

The function `write.table` provides a simple interface for writing data to a file. It can be used to write a dataframe or a matrix to a file. For example, if `mydats` is a dataframe, the command

```

> write.table(mydats, file="my_data.txt",
+            col.names=TRUE, row.names=FALSE)

```

will store it in the text file `my_data.txt` with the labels for the variables or factors it contains in the first row (`col.names=T`), but without the numbers associated with each row (`row.names=F`). The `sep` option is used to choose the separator for the data, the default is a blank space `sep=" "`, but you can choose a comma `sep=","` a semicolon `sep=";"` or other meaningful separators.

5.2.2 Low-Level File Output

Perhaps the most user friendly low-level function for writing to a file is `cat`. Suppose you have two vectors, one with the heights of 5 individuals, and one with an identifier for each, and you want to write these data to a file:

```
> height <- c(176, 180, 159, 156, 183)
> id <- c("s1", "s2", "s3", "s4", "s5")
```

you can use `cat` in a for loop to write the data to the file, but let's first write an header

```
> cat("id height \n", file="foo.txt", sep=" ", append=FALSE)
```

the first argument is the object to write, in this case a character string with the names of our variables to make a header, and a newline character to start a new line. Now the for loop:

```
> for (i in 1:length(id)){
+   cat(id[i], height[i], "\n", file="foo.txt",
+       sep=" ", append=TRUE)}
```

notice that this time we've set `append=TRUE` to avoid overwriting both the header, and any previous output from the preceding cycle in the for loop. We've been using a blank space as a separator, but we could have used something else, for example a comma (`sep=','`).

Chapter 6

Descriptive

6.1 Tables

6.2 The `scale()` Function

The `scale()` function can be used to easily transform your data into z scores. Here's a silly example

```
> a<-c(1,2,3)
> scale(a)
      [,1]
[1,]  -1
[2,]   0
[3,]   1
attr(,"scaled:center")
[1] 2
attr(,"scaled:scale")
[1] 1
```

Chapter 7

Graphics

7.1 Overview of R Graphics

Table 7.1: R Graphics Functions

Function
<code>plot()</code>
<code>barplot</code>
<code>boxplot</code>
<code>histogram</code>
<code>matplot</code>
<code>stripchart</code>
<code>interaction.plot</code>

7.2 The `plot()` Function

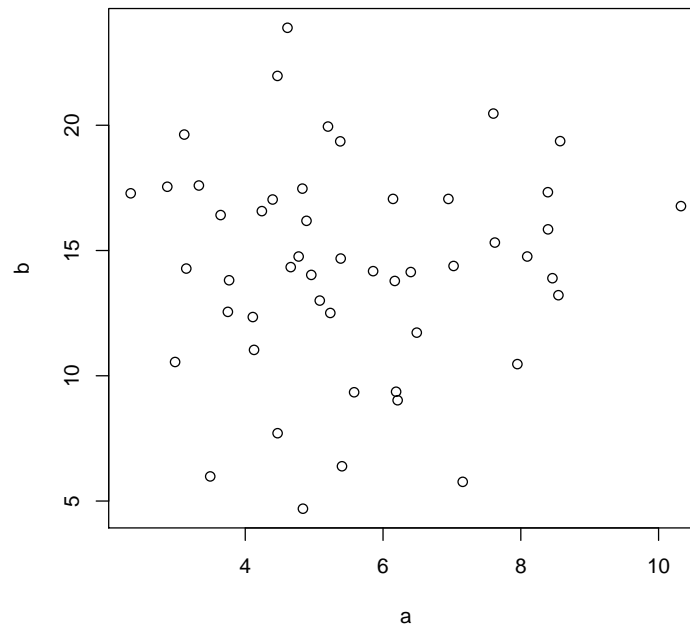
The `plot()` function is most commonly used to draw a scatterplot of two variables, however if given an R object with a plot method as an argument it will produce different types of graphics depending on the object it is plotting. Let's see an example of a scatterplot with some fake data:

```
a <- rnorm(5,1.6,n=50)
b <- rnorm(15,4.3,n=50)
```

this creates two vectors of length 50 with values normally distributed, now we can plot the values of vector `b`, against the values of vector `a` with:

```
> plot(x=a,y=b)
> ## or for short:
> plot(a,b)
```

Figure 7.1: A Scatterplot



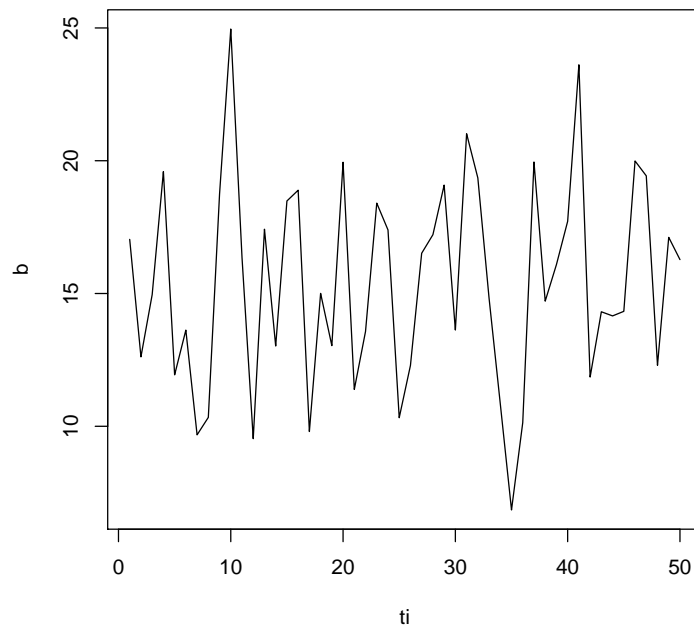
the resulting scatterplot appears in Figure 7.1

If we were to plot the change of a variable over time, it could be a good idea to connect the values at different time points in the plot with lines, this is easily achieved setting the option `type`. Below is an example, the result is shown in Figure 7.2.

```
> ti<-1:50  
> b<-rnorm(15,4.3,n=50)  
> plot(ti, b, type='l')
```

Some other possible values for the option `type` are 'p' for points (the default), 'b' for both points and lines, and 'o' for overplotted points and lines (very similar to 'b').

Figure 7.2: Values Connected by Lines



7.3 Drawing Functions

The `plot` function can be used for drawing mathematical functions, for example:

```
> vec <- seq(from=0, to=4*pi, length=120)
> plot(vec,sin(vec), type="l")
```

in this case is necessary to use 'l' as the type of plot, otherwise the plot would resemble a messy scatterplot. It is also possible to plot markers on the points in which the function is actually evaluated. Using the option 'b' for the plot type, both a continuous line and markers are plotted. There are lots of options to define the appearance of a plot, next comes an overloaded example used just to introduce some of these options. A more detailed description is given in Section 7.9.

7.3.1 `matplot`

Plotting the sine and cosine functions together with the `matplot` function, the result is in Figure 7.3

```
> a<- seq(from=0, to= 2*pi, length=20)
> s <- sin(a)
> c <- cos(a)
> aa <- cbind(a,a) #matrix of hor coord
> cs <- cbind(c,s) #matrix of vert coord
> matplot(x=aa,y=cs,type="l",lwd=1.8,ylab="sine and cosine function")
```

7.4 Barplots

Barplots sometimes come in handy when you want to summarise your data. Suppose we have administered a test to three different groups of people, which we will designate as "a", "b" and "c". We want to summarise and compare the performance of each group with a nice graph, a barplot will make it. The data are stored in the file "test.txt", in the format shown in table 7.2

We can read in the file directly as a dataframe:

```
> test <- read.table("test.txt", header=TRUE)
```

We want to use the `tapply` function to get quickly summary tables with the means and standard deviations for the three groups. However the format of the data frame at this point is not suitable for the `tapply` function, because it has 3 observations for each row, and the `tapply` function can be used only

Figure 7.3: sine and cosine functions with matplotlib

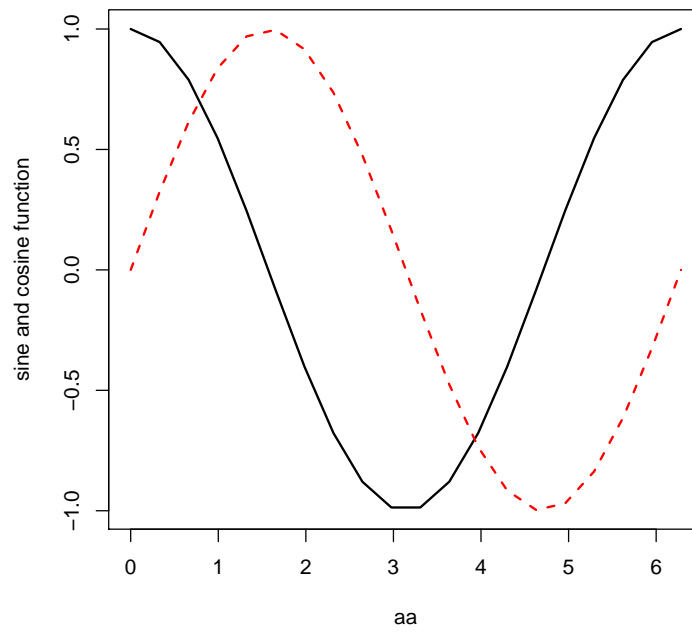


Table 7.2: Data for the test example

a	b	c
4	6	5
5	8	3
3	7	8
6	5	4
5	9	9
7	7	8
5	6	5
8	5	7
5	8	6
4	10	4

with the format “one row per observation”, in which we have the values of the observations in one column and a set of “labels” identifying the group to which a given observation belongs to in another column. Fortunately, we can easily change the format of our dataframe with the `stack` command. What it does is just to create a single “values” vector from the three columns we had previously, and to add automatically another “index” vector with the labels we need:

```
> test2 <- stack(test)
```

we can have a look at the new dataframe simply typing:

```
> test2
```

please, notice that the `stack` function has automatically named the two vectors `values` and `ind`, we need to know these names to use the `tapply` function.

Now we’ll create the two summary tables using the `tapply` function, one with the means and one with the standard deviations for the three groups:

```
> test_means <- tapply(X=test2$values, IND=test2$ind, FUN=mean)
> test_sd <- tapply(X=test2$values, IND=test2$ind, FUN=sd)
```

Now we can draw a simple barplot displaying the means for each group:

```
> barplot(test_means, col=c("darkred","salmon2","plum4"))
```

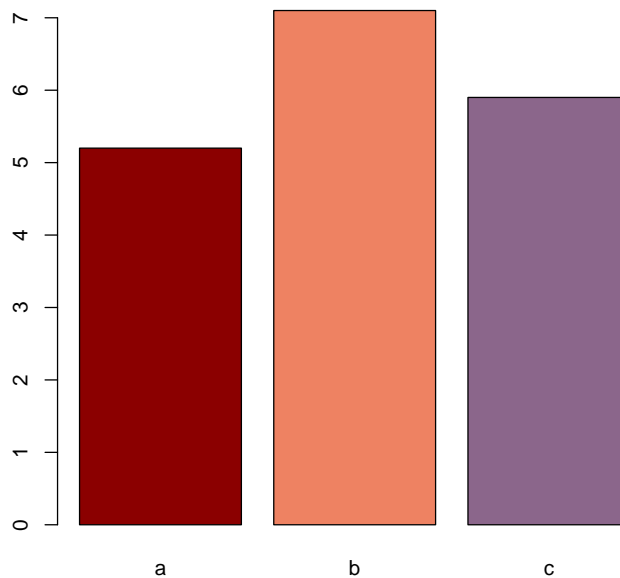
we’ve added colours to the graph using the `col` option. You can look at the resulting graph in figure 7.4 You can control the width of the bars specifying the `width` option, and setting the range for the x axis with the `xlim` option, specifying only the width does nothing, you must also set the `xlim`. You set `xlim` with a vector of the form `xlim = c(from, to)`, in which “from” is the origin and “to” is the end of the axis. In the example below we will set `xlim` to go from 0 to 3 and the bars to have a width of 0.5:

```
> barplot(test_means,col=c("darkred","salmon2","plum4"),
+         xlim = C(0,3), width=0.5)
```

this sets the width of all bars at 0.5, we could specify the width for each single bar instead, by giving to ‘width =’ a vector with the width values for each bar.

You can also control the spacing between the bars with the `space` option. The default is set to 0.2.

Figure 7.4: Simple Barplot



Barplots with Error Bars

Now let's say we want to get the same barplot but with error bars showing the standard deviation for each group. We could achieve this result adding lines to the current graph, but there is a better option, we can use the `barplot2` command, which comes with the `gplots` library and provides an easy way of adding error bars to a barplot. So once we have the `gplots` package installed we first load it:

```
> library(gplots)
```

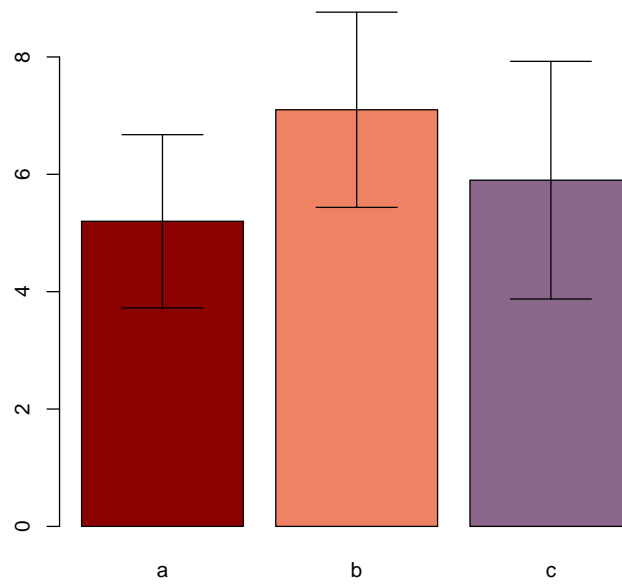
and then we can draw our barplot with error bars. To get them, we need to set the option `plot.ci=TRUE` and then specify the upper and lower bounds of the error bars with the `ci.u` and `ci.l` commands. So we first create the values for `ci.u` and `ci.l`, so that the error bars represent one standard deviation around the mean. We'll use the values from the two tables we had created before, with means and standard deviations for the three groups:

```
> upper <- test_means + (test_sd)
> lower <- test_means - (test_sd)
```

Now we can draw our barplot, you can see the result in figure 7.5:

```
> barplot2(test_means,col=c("darkred","salmon2","plum4"),  
+          plot.ci=TRUE,ci.u=upper,ci.l=lower)
```

Figure 7.5: Barplot with Error Bars



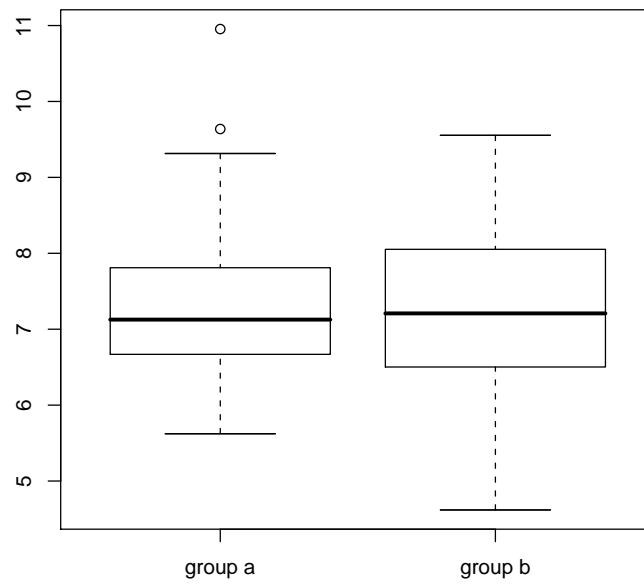
7.5 Boxplots

Boxplots provide a powerful way to visualise the central tendency and the dispersion of the data for a given sample, and to directly compare these same characteristics for different samples. Let's look at one of them, the data are organised in a dataframe in the file `boxplot1.txt`. They are the scores for two different groups (group a and group b) in a test. With the boxplot we want to see how the scores are distributed in the two groups. The dataframe contains a column `score`, with the score for each subject and another column `group` that defines the group a given observation comes from. So we first read in the dataframe, and then ask for the boxplots with the distribution of scores, as a function of the group they belong to.

```
> dats <- read.table("boxplot1.txt", header=TRUE)
> boxplot(dats$score~dats$group, names=c("group a", "group b"))
```

The results are in figure 7.6, the thick black lines in the middle of the boxes represent the median, if this is about the middle of the box, the distribution of the data should be normal. The two lines that delimit the box are called “hinges”, and they are approximately the first and the third quartiles. The horizontal lines that form the ‘Ts’ above and below the box are called “whiskers”, and inside them are contained all the observations that fall within a distance of 1.5 times the size of the box, upwards or downwards. Points that fall outside this distance are outliers and they are represented as a circle. In our case there are two outliers in group a. Apart from checking if the distribution is normal, you can also check if the variances are approximately equal, by comparing the size of the boxes (the distance between the hinges). If one boxplot is clearly bigger than the other one (for example two times bigger), then the variances for the two groups are likely not to be equal.

Figure 7.6: Boxplots comparing the distribution for two groups

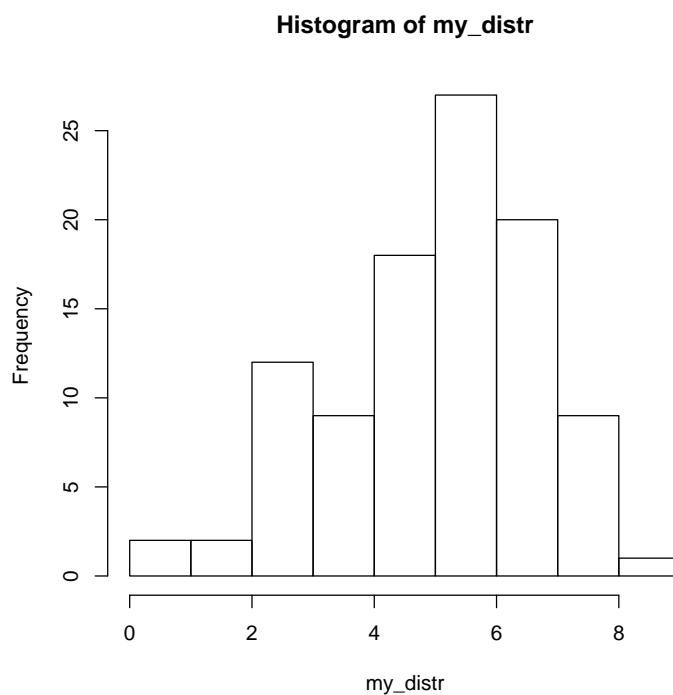


7.6 Histograms

Histograms are a good way to visualise the distribution of a sample, the function `hist()`, can be used to plots a histogram of frequencies (counts) of the sample, or of its density function (setting the option `freq=F`). Let's first create a sample with a normal distribution, and then plot its histogram, the result is in Figure 7.7.

```
> my_distr <- rnorm(100, 5, 1.7)
> hist(my_distr)
```

Figure 7.7: Frequency Distribution of a Random Sample



7.7 Stripcharts

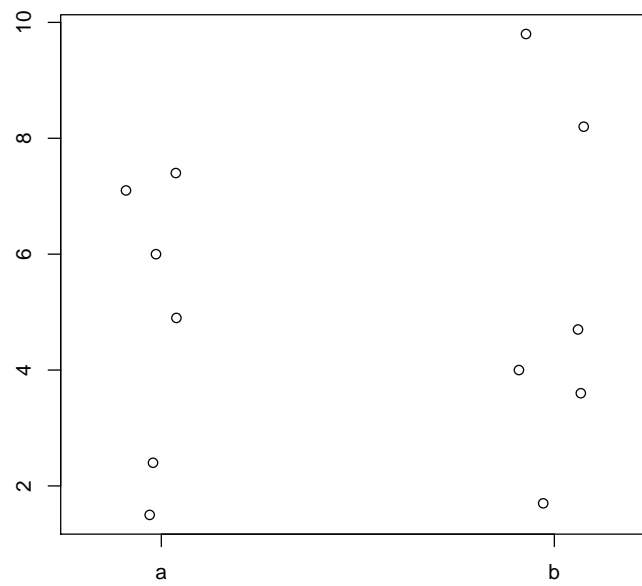
If the groups contain a small number of observations, it might be better to use a stripchart to visualise their distributions. In a stripchart each point represents a single observation. By default they are drawn on a line, so if two observations have the same score, they overlap. To avoid this problem you can give a certain amount of jitter to the plot, so that observations with the same score are scattered a little and can be easily distinguished. Here's the code for producing a stripchart, the data are in the file `stripchart1.txt`

and they are arranged in the same way as the data in `boxplot1.txt`, just the sample sizes are smaller, with 6 observations per group.

```
> stripchart(dats$score~dats$group, method="jitter",  
+           jitter=0.1, pch=1, vertical=TRUE)
```

The resulting graph is shown in figure 7.8

Figure 7.8: Stripchart 1

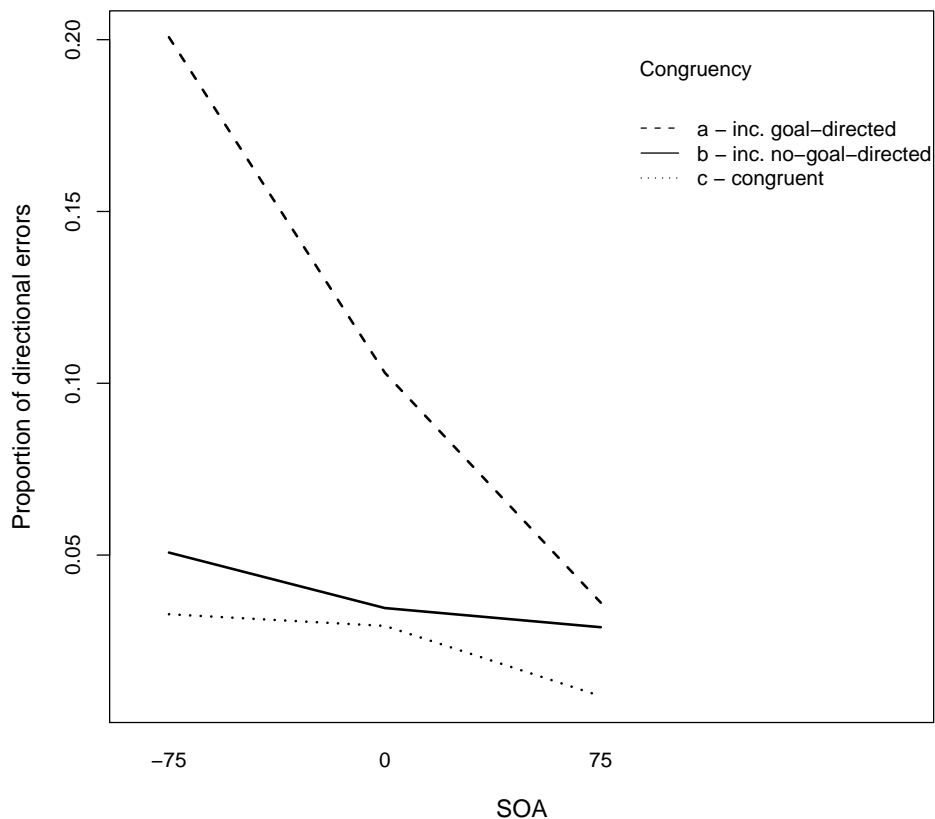


7.8 Interaction Plots

Interaction plots can be used to visualise the means for the levels of a factor, at the levels of another factor, for example in a two-way ANOVA design, and in this way they allow to spot possible interactions between the two factors involved. The following graph comes from a two-way ANOVA design, in which the variable of interest was the proportion of errors in a task, as a function of the spatial congruency and SOA (stimulus onset asynchrony), of a distracting stimulus presented during the execution of the task.

```
> interaction.plot(dats$soa,dats$congr,dats$errors,  
+                 ylab="Proportion of errors", xlab="SOA",  
+                 trace.label="Congruency",  
+                 lty=c("solid","dashed","dotdash"))
```

Figure 7.9: Interaction Plot with Different Line Types

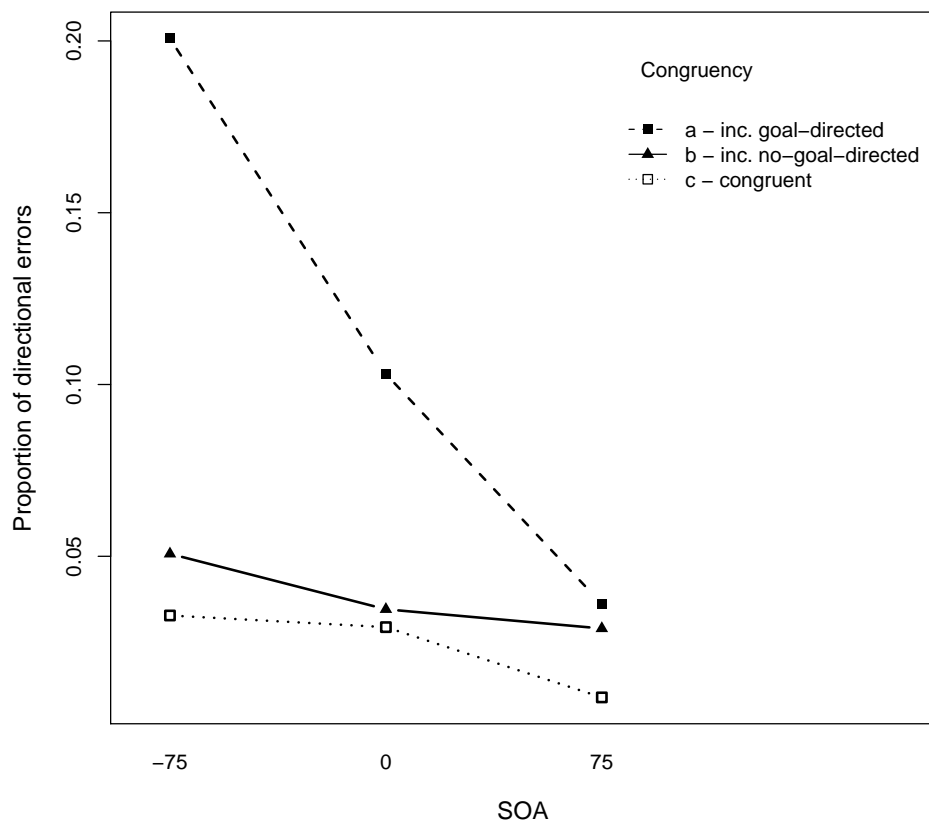


well the essential ingredients to get the plot are just the first three arguments (the others are optional), the default order in which you have to give them

is a bit awkward, because the variable of interest (in our case `dat$errors`) is the third argument, the first one is the factor that goes on the x axis, and the second one is the *trace factor*, whose levels will be represented as lines of a different type, or of different colours. The y axis yields the measure for our variable of interest. You can see the plot in Figure 7.9. Another good way to represent the levels for the trace factor, is through the use of symbols, in the following graph (see the result in Figure 7.10) both line type `lty` and symbols `pch` are used to differentiate between the levels of the trace factor, in order to get this you need to set the option `type` to `"b"` that means use both line type and points (symbols), setting this option to `"l"` will give just different line types while setting it to `"p"`, will give just different symbols. In these examples I specified the line types and the symbols to use, but this is not necessary, if you don't, R will cycle through its different line types and/or symbols as necessary to represent all the levels of the trace factor. See Section 7.9.1 for a description of the different line types and Section 7.9.2 for a description of the different symbols available in R.

```
> interaction.plot(dats$soa,dats$congr,dats$errors,  
+                 ylab="Proportion of errors", xlab="SOA",  
+                 trace.label="Congruency",  
+                 lty=c("solid","dashed","dotdash"),  
+                 type="b",pch=c(0,15,17))
```

Figure 7.10: Interaction plot with different Line Types and Different Symbols



7.9 Setting Graphics Parameters

Graphics parameters allow you to tweak many elements of a plot, such as the font for the labels, the symbols or line types to use and so on, see `?par` to get a full list and description of these parameters. Graphics parameters can be set and accessed with the function `par`, called without arguments, as `par()`, it will give you a full list of the current defaults, if you want to query only one or a few parameters use:

```
> par("lwd")          ## see current line width
> par(c("lwd","pch")) ## see lwd and plotting symbols
```

to change the value of a parameter you can use:

```
> par(lwd=1.4)        ## change line width
> par(pch="*",        ## change plotting symbol
+     bg="gray80")    ## use a light gray background
```

Moreover most plotting functions like `plot`, `barplot` and so on, allow you to set some of the parameters for the current plot, as an argument to the function itself, for example in `plot` you can choose the type of plot (points vs lines) and the plotting symbol as options with `type` and `pch`:

```
> d <- rnorm(15,4,2)
> e <- rnorm(15,9,3)
> plot(d~e, type="p", pch=3)
```

The following sections will give a more in depth explanation of some graphics parameters, before that however we'll have a closer look at how to use `par`.

Saving and Restoring Graphics Parameters

Often you'll want to change the graphics parameters only for a few plots, and then reset them back to the defaults. The function `par` when used to change the value of some graphics parameters returns a list with the *old* values of the graphics parameters that have changed:

```
> par("lwd", "col") #these are the default parameters
$lwd
[1] 1

$col
[1] "black"
> oldpar <- par(lwd=2, col="red") #while changing the parameters
>                                     #we store the old values in a list
> oldpar
```

```

$lwd
[1] 1

$col
[1] "black"
> s <- seq(0, 10, .1)
> plot(s,sin(s))           #we plot something
> par(oldpar)             # and then we restore
>                          # the old parameters

```

notice that calling `par`, opens a graphics device if there is not one already open, the changes you do using `par` apply only to this graphics device, any other new graphic device that you open will have the default graphics parameters.

List of Graphics Parameters by Category

Table 7.4: Parameters for Colours

Parameter	Function
<code>col</code>	plotting colour
<code>col.axis</code>	colour for axis annotation
<code>col.lab</code>	colour for x and y labels
<code>col.main</code>	colour for main title
<code>col.sub</code>	colour for sub-titles
<code>bg</code>	background colour
<code>fg</code>	foreground colour

Table 7.5: Parameters for Fonts

Parameter	Function
<code>family</code>	font family
<code>font</code>	type of font (plain, bold, italic etc...)

7.9.1 Line type `lty`

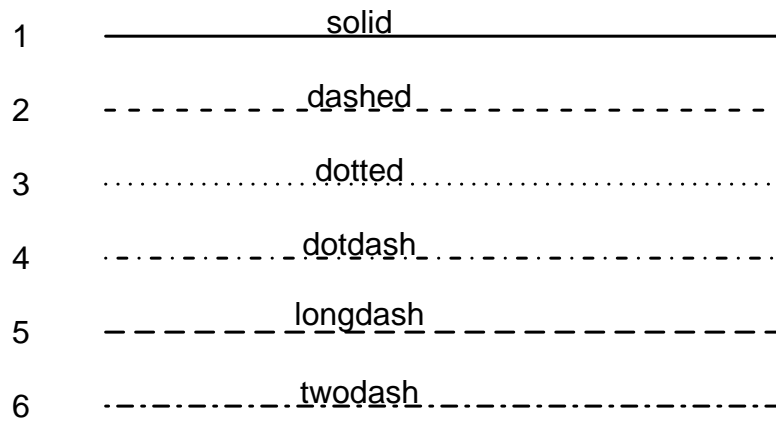
There are six line types that you can call in R just with names or numbers¹. These are listed in Table 7.6 and shown in Figure 7.11. There is a different, more complicated way for setting many more line types, please, consult the manual for further information on that.

¹Actually there are seven, the first one is in fact “blank” or 0 which just draws nothing.

Table 7.6: The Six Default Line Types in R

N.	Name
1	solid
2	dashed
3	dotted
4	dotdash
5	longdash
6	twodash

Figure 7.11: The Six Default Line Types in R



7.9.2 Symbols with pch

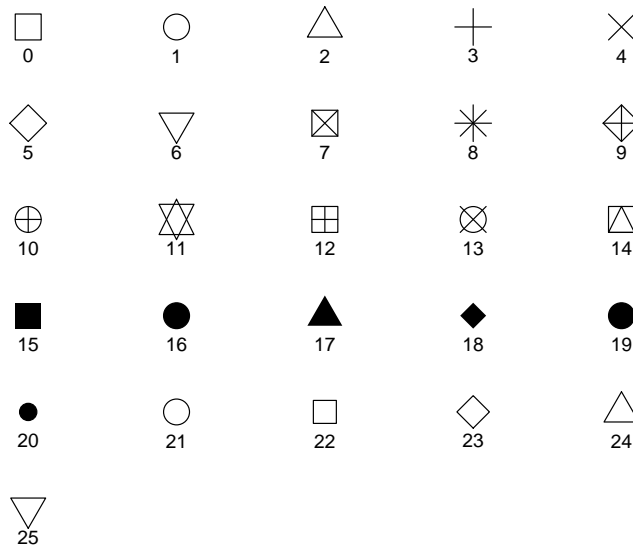
`pch` is a graphical parameter for changing the way points are plotted in certain graphical functions. This parameter can be set in two ways, the first one is to give a symbol to be plotted as a character, for example

```
pch="+"  
pch="T"  
pch="*"  
pch="3"
```

in this case you enclose the character you want to use between quotes, you can't use more than a single character. The other way to set the `pch` is to use a number between 0 and 25, which will select one of 26 special symbols available for plotting (see Figure 7.12), like circles, triangles, and so on:

```
pch=1  
pch=3  
pch=5
```

Figure 7.12: Plotting Symbols in R



7.10 Adding Elements to a Plot

7.10.1 Adding a Legend

Some graphics functions by default add a legend to the graph (e.g. `interaction.plot`), or allow to add (or discard) a legend by setting an option inside the function (e.g. `barplot`). However the default settings for the legend, such as positioning, text or symbols, might not be suitable to your graph, in which case, you need to turn off the default legend (if there is one), and add a customised legend with the `legend` function. Below there is an example of an interaction plot with the legend added through the `legend` function, the resulting graph is in Figure 7.13

```
>interaction.plot(dats$isi,dats$congr,dats$pcdirerr,  
+ ylab='Mean proportion of directional errors', xlab='SOA',  
+ type='b',pch=c(0,15,17),legend=FALSE) ## eliminate default legend  
>legend('topright',legend=c('a - congruent','b - inc. goal-directed',  
+ 'c - inc. no-goal-directed'),lty=c(3,2,1), pch=c(0,15,17),  
+bty='n', title='Congruency')
```

here we first made the plot, and afterwards we added the legend. The first argument, `topright` indicates the position where we want the legend to appear, other possible values are `bottomright`, `bottomleft`, `right`, `bottom`, `center` and so on. It is also possible to specify the position of the legend by giving the coordinates of its top-left corner, for the above example we might have written:

```
> legend(x=2.3, y=0.2, legend=c('a - congruent','b - inc. goal-directed',  
+ ....example not continued
```

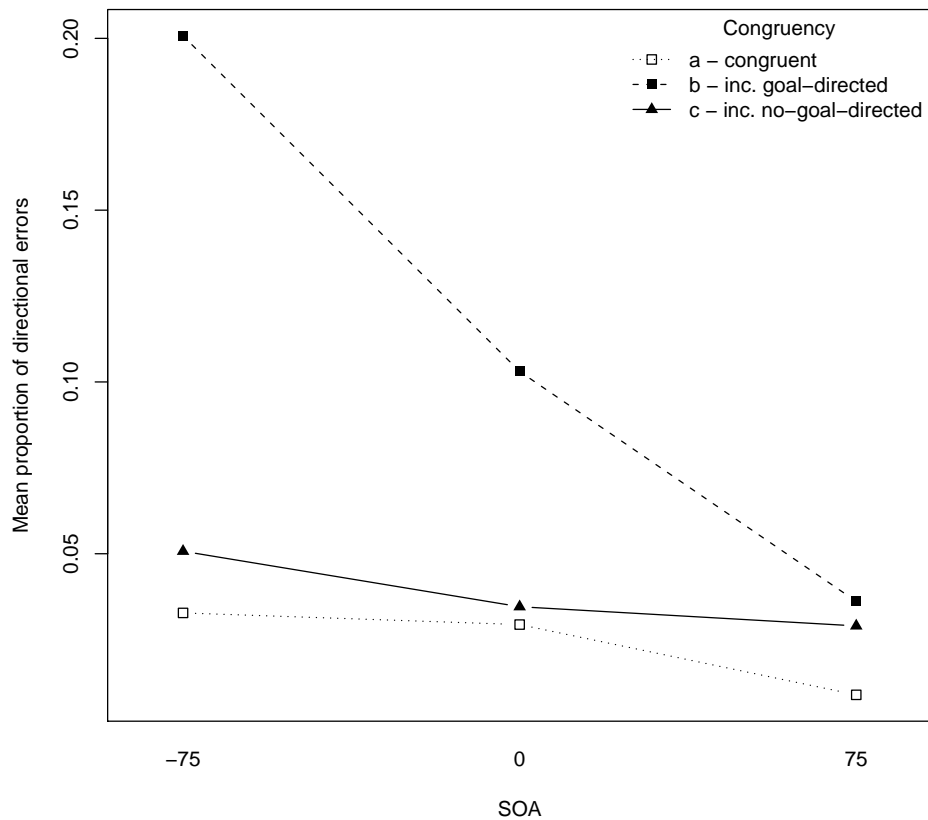
The text of the legend is given as a character vector, each element of the vector represents one item of the legend.

Since we're using both different line types, and different symbols to differentiate the lines in the interaction plot, for the legend we need to specify both in the legend, with the `lty`, and `pch` arguments, the line types and symbols of course, should be the same as those used for the plot.

In the above example we suppressed the drawing of a box around the legend setting the `bty` argument to `n`, if you want a box set it to `o` (this is the default anyway, so you don't really need to specify it if you want the box). If you choose to enclose the legend into a box, you can set its background colour through the `bg` argument.

There are of course many other options, to get a full listing, please, look up the manual.

Figure 7.13: Plot with "manually" added legend



7.10.2 Adding Text

You can insert text in a graph with the `text` function, you have just to specify the x and y coordinates of the point on which to center the text:

```
> plot(x,y)
> text(x=3,y=1.5, "mean for control group")
```

if you want to use \LaTeX mathematical symbols, you have to use the `expression` function:

```
> text(x=3,y=1.5,expression(alpha))
```

7.10.3 Adding a Grid

A grid can be easily added to an existing plot with the function `grid()`

```

> s<-seq(from=0, to=2*pi, length=100)
> plot(s, sin(s))
> plot(s, sin(s), type='l')
> grid() ##add the grid

```

The default colour for the grid is lightgray, you can choose another colour setting the 'col' option:

```

> grid(col='red')

```

7.10.4 Setting the Axes

If you don't like the way the axes are set for a given plot, you can draw the plot without them first, and then add customised axes with the `axis` function. There are several ways to get rid of the default axes on a plot:

```

> plot(1:10, axes=FALSE) #do not draw any axes or box around
>                               # the plot
> plot(1:10, xaxt="n")      #don't draw the x axis alone
> plot(1:10, yaxt="n")      #don't draw the y axis alone

```

Once you've got rid of one or more axis, you can draw them calling `axis`:

```

> axis(1, at=seq(1, 10, 3), labels=as.character(seq(1, 10, 3)))

```

the first argument specifies the side on which the axis should be drawn, 1 means the bottom axis, 2 the left axis, 3 the top axis, and 4 the right axis. See `?axis` for other arguments to the function.

7.11 Creating Layouts for Multiple Graphs

7.11.1 mfrow and mfcop

The parameters `mfrow` and `mfcop` allow you to divide the graphics device you're using (e.g. X11 or pdf) into multiple boxes, each box will contain a new figure. These parameters are set giving a vector of the form:

```

> par(mfrow=c(n_rows,n_columns))

```

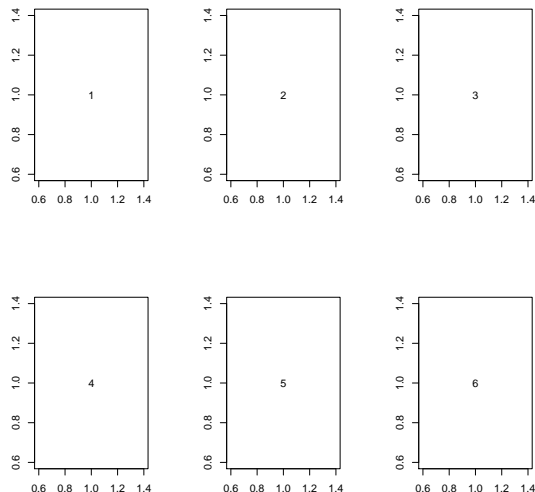
where `n_rows` is the number of rows and `n_columns` is the number of columns you want for your layout. For example, the following will create a layout with 2 rows and 3 columns as you can see in Figure 7.14:

```

> par(mfrow=c(2,3))
> symb <- as.character(1:6)
> for(i in 1:6){
  plot(1, 1, pch=symb[i])
}

```

Figure 7.14: A 2x3 Layout with `mfrow`



`mfcol` is exactly the same as `mfrow` but the figures are drawn in sequence by column rather than by row, for example the following code yields the Figure 7.15

```
> par(mfcol=c(2,3))
> symb <- as.character(1:6)
> for(i in 1:6){
  plot(1, 1, pch=symb[i])
}
```

of course you can get other layouts, try:

```
> par(mfrow=c(2,2)) ## 4 nice boxes
> par(mfrow=c(1,3)) ## one row, 3 cols
```

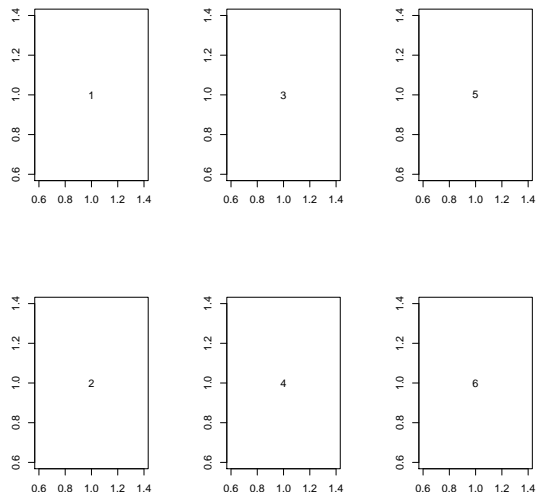
Rather than having the figures drawn sequentially, following the order determined by `mfrow` or `mfcol`, it is possible to specify directly the position for the next figure using the `mfg` parameter. The next example will draw the plot in the slot defined by the crossing between the second row and the first column of a 2 x 2 layout:

```
> par(mfrow=c(2,2)) ##create 2x2 layout
> par(mfg=c(2,1)) ##set next fig at 2dn row, 1st col
> plot(1:100, sin(1:100))
```

7.11.2 layout

A more flexible way to divide a graphics device into multiple plotting regions is given by the `layout` function. With the `layout` function, the graphics

Figure 7.15: A 2x3 Layout with `mfc`



device is divided into a matrix of n_rows x n_cols sub-regions, and each figure is assigned to one or more of these sub-regions. Let's look at an example

```
> m<-matrix(c(1,2,3,4),nrow=2,byrow=TRUE)
> m
      [,1] [,2]
[1,]    1    2
[2,]    3    4
> layout(m)
```

the matrix `m` we've created, completely defines our layout for the graphics window. The window is divided into $2 \times 2 = 4$ sub-regions. Moreover, the first figure is assigned the sub-region on the top-left of the window, the second figure the sub-region at the top-right, the third the region at the bottom-left, and the fourth the region at the bottom-right. This example is in itself not very different from what we would get with `mfrow`, however two key differences make `layout` more powerful than `mfrow`. First, it is possible to assign more than a single sub-region to a figure, for example

```
> m <- matrix(c(1,1,2,3), nrow=2, byrow=TRUE)
> m
      [,1] [,2]
[1,]    1    1
[2,]    2    3
> layout(m)
```

divides the graphics window into 4 sub-regions as above, but the first figure is assigned the two sub-regions on top, while the bottom-left sub-region is assigned to the second figure, and the bottom-right to the third.

The second key difference with `mfrow` is that with `layout` it is possible to define the width of the columns, and the height of the rows composing the array of sub-regions in the graphics window. Width and height are given as vectors of *relative* widths and heights. For example:

```
> m <- matrix(c(1,1,2,3), nrow=2, byrow=TRUE)
> m
      [,1] [,2]
[1,]    1    1
[2,]    2    3
> layout(m, width=c(1/4,3/4), height=c(2/3,1/3))
```

will make the first column 1/4 of the total width, and the second column the other 3/4, the same reasoning applies to the row heights. The command `layout.show(n)`, where n is the number of n^{th} figure that is going to be plotted, will show the outline of its layout in the graphics window.

7.12 Graphics Device Regions and Coordinates

In traditional R graphics the graphics device (e.g. the `X11()` window where your plot and annotations appear) is divided into different regions (see Fig. 7.16:

- the *plotting region* is the area in which the drawing of points or lines representing your data occur. The plotting region is contained into the *figure region*
- the *figure region* is composed of the plotting regions plus the margins where the plot can be annotated with labels for the axes, a title etc. . .
- the *outer margins* surround the figure regions. The outer margins are usually set to zero (i.e. there are no outer margins), they become useful however for annotating multiple plots that appear on the same page (e.g. plots generated with `mfrow`). When multiple plots are arranged on the same page (or device), each is assigned a *figure region* with its own margins, so the outer margins can be used for annotating the overall page (see Fig. 7.17).

The width and height of the device is usually specified when the device is opened, for example

```
> X11(width=8, height=8)
```

opens a `X11()` device measuring 8x8 inches. The size of an open device can be queried with `par("din")`, this is a read only graphics parameter, which I guess means that once a certain device is opened its size cannot be changed (an `X11()` window however can be re-sized with the mouse, and `par("din")` correctly reports the new size). Different units of measure can be used to specify the size of the areas inside a device, some of these units will be shortly introduced here, others will be explained when they are first used:

- *inches*: an inch is 2.54 centimetres (notice that the actual physical measure of what you see on your monitor may depend on your monitor's settings, e.g. dpi, screen resolution, etc...)
- *lines of text*: this measure depends on the value of `cex` and `pointsize`
- *Normalised Device Coordinates (NDC)*: the device region is 1x1 NDC whatever the actual physical measure. The lower left corner has coordinates (x=0, y=0) and the upper right corner (x=1, y=1). Using NDC thus the size of regions inside the device can be specified in relative terms to the device size.

In the next paragraphs the graphics parameters for controlling the different regions inside the device will be explained. Always keep in mind the layout of a graphics device in R (see Fig. 7.16 and Fig. 7.17).

Figure Region

The figure region can be set either in inches or in normalised device coordinates.

- **fig**: NDC coordinates of the device in the form `c(x1, x2, y1, y2)`, where (x1, y1) are the coordinates of the lower left corner, and (x2, y2) are the coordinates of the upper right corner). Example:

```
> par(fig=c(0.1, 0.5, 0.1, 0.6))
> plot(1:10)
> par(fig=c(0.5, 1, 0.1, 0.6)) # the next call to plot will erase the
>                               # current plot
> plot(1:10)
> par(fig=c(0, 0.5, 0.1, 0.6), new=TRUE) # the next call to plot will not
>                                           # erase the current plot
> plot(1:10)
```

as shown in the example to change the figure region without starting a new plot add `new=TRUE`, this may be used for creating complex arrangements for multiple plots within the same device.

- `fin`: the figure region dimension (width, height) in inches. Example:

```
> par(fin=c(5,5))
```

Plotting Region

- `plt`: a vector of the form `c(x1, x2, y1, y2)` giving the coordinates of the plot region as fractions of the current figure region
- `pin`: the current plot dimensions, (width,height), in inches

Margins

- `mar`: the width of the margins for the four sides of the plot, specified in terms of *lines of text*. The margins are specified in the form `c(bottom, left, top, right)`. The default is `c(5,4,4,2) + 0.1`. Example:

```
> par("mar") # get current margins size
[1] 5.1 4.1 4.1 2.1
> par(mar=c(6, 2, 3, 0) + 0.1) #set new margins size
```

- `mai`: the same as `mar`, but the unit of measure is inches rather than lines of text

Outer Margins

- `oma`: a vector of the form `c(bottom, left, top, right)` giving the size of the outer margins in lines of text
- A vector of the form `c(bottom, left, top, right)` giving the size of the outer margins in inches
- A vector of the form `c(x1, x2, y1, y2)` giving the outer margin region in normalised device coordinates (NDC)

Figure 7.16: The Figure Region Includes the Plotting Region and the Margins

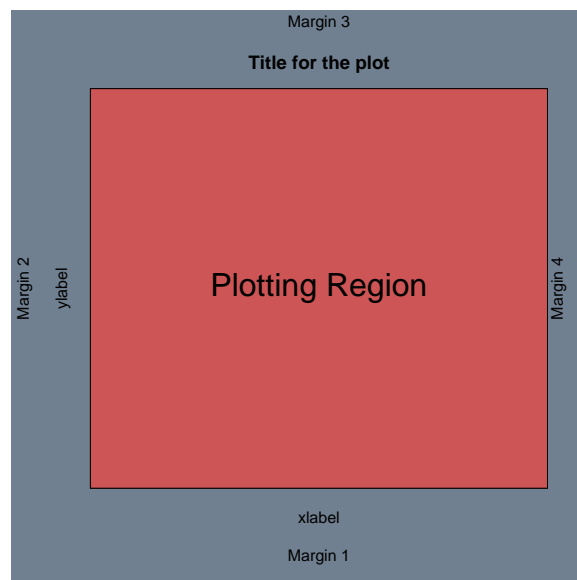
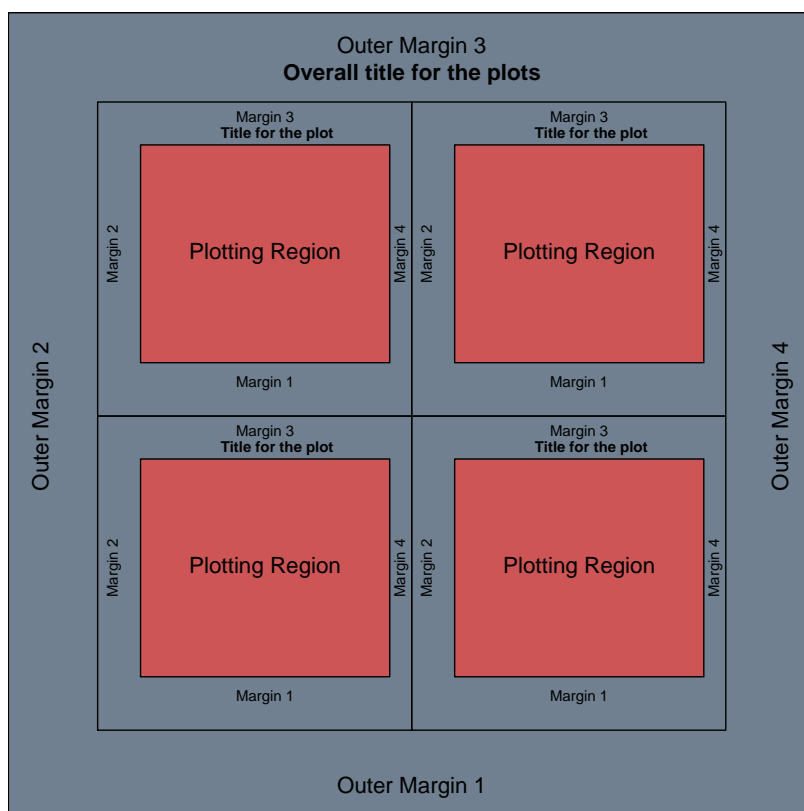


Figure 7.17: Device with Outer Margins and Multiple Figure Regions



7.13 Plotting from Scratch

The high level plotting functions such as `plot`, `histogram`, `barplot` and so on, provide a good and quick way to produce graphs. Plotting from scratch, using the low-level plotting commands is generally not necessary, unless you want to create some new, customised plotting functions. Learning to plot “from scratch”, however is a very good way to learn how graphics parameters work, which is often necessary to customise plots created with the high-level plotting functions.

We’ll start with a very simple example of a scatterplot

```
> plot.new()
> plot.window(xlim=c(0,10), ylim=c(0,10))
```

the `plot.new()` command creates a frame for plotting, and opens a graphics device if there is not one already opened. `plot.window` defines the limits for the x and y axes, points outside these limits will not appear in the plot. After these two commands we’re ready to do the actual drawing

```
> points(x=c(1,2,3,4,5,9), y=c(2,5,3,4,5,3))
> axis(side=1)
> axis(side=2)
```

`points` will draw points at the coordinates given in the `x` and `y` arguments. To complete this very minimal plot you need at least some axes. The `axis` function adds the axis, the `side` argument specifies where the axis should be drawn, 1 means at the “bottom”, 2 at the “left” side, and so on in a clockwise fashion.

7.14 Colours for Graphics

The command `colours()` gives a list of built-in colours available for graphics in R. You can see some of these colours in Table 7.7. There are 101 built-in shades of gray, from `gray0`, that is almost black, to `gray100` that is almost white, you can see some of them in Table 7.8. A complete table of built-in R colours is given in Appendix E.

You can also specify colours in `rgb` values. By default R accepts values in the range 0–1, but you can change the range with the `max` option to set the range as 0–255. Please note that changing the range doesn't change the colours you can use, it just changes the values you use to specify them, so for example the following graphs will have the same colours:

```
> vec <- c(3,6)
> barplot (vec, col= c(rgb(0.176, 0.262, 0.49),
+                       rgb(0.568, 0.254 0.654)))
> barplot (vec, col= c(rgb(45,67,125, max=255),
+                       rgb(145,65,167, max=255)))
```

the first uses the default range, and the second uses the range 0–255, but I simply derived the values for the first graph, dividing those for the second by 255.

The function `col2rgb` can be used to get the `rgb` values of a built-in colour from its name. The `rgb` values are given in this case in the range 0–255. Here's an example:

```
> col2rgb("lightslateblue")
      [,1]
red      132
green    112
blue     255
```

Table 7.7: Some built-in colours in R

 "aliceblue"	 "honeydew2"	 "olivedrab"
 "antiquewhite2"	 "indianred"	 "orange"
 "aquamarine"	 "ivory"	 "papayawhip"
 "beige"	 "khaki2"	 "peachpuff"
 "bisque"	 "lightgreen"	 "peru"
 "brown3"	 "lightgoldenrod"	 "plum4"
 "burlywood2"	 "lavender"	 "powderblue"
 "chocolate4"	 "lemonchiffon"	 "rosybrown"
 "cornflowerblue"	 "lightblue"	 "salmon2"
 "cornsilk"	 "lightcyan"	 "sandybrown"
 "coral2"	 "linen"	 "seashell"
 "darkkhaki"	 "lightslateblue"	 "sienna"
 "darkolivegreen"	 "magenta"	 "slategray"
 "darkred"	 "mediumaquamarine"	 "tan1"
 "dodgerblue"	 "mediumorchid"	 "thistle1"
 "firebrick2"	 "mediumpurple"	 "tomato2"
 "forestgreen"	 "midnightblue"	 "turquoise"
 "gold"	 "mistyrose"	 "wheat"

Table 7.8: Different shades of gray

 "black"	 "gray40"	 "gray75"
 "gray9"	 "gray50"	 "gray85"
 "gray25"	 "gray60"	 "gray94"
 "gray30"	 "gray65"	"white"

7.15 Managing Graphic Devices

7.15.1 Opening Another Graphics Window

When you issue the command for a graph R opens a window to show it, if you afterwards issue another command for a graph, if the previous window is still open, R doesn't open another one, but rather replaces the old graph with the new one. If you wish to show the new graph in a separate window, you have to open the graphic device yourself, this is accomplished with the command `X11()` under Unix and with the command `windows()` under the Windows OS. The device window can also be closed from the command line with

```
> dev.off()
```

if you have many device windows open and you want to close them all at once use

```
> graphics.off()
```

For further functions to manage multiple device windows see `?dev.set`.

7.15.2 Exporting Graphics

With R it's also possible to export your graphics in different file formats, such as JPEG or postscript files. To do this, you need to open first the graphics device you want to use, then insert the command for the graphic, and finally turn off the graphic device. Here's an example of how to produce a graphic in JPEG format:

```
> jpeg(file="plot.jpeg")
> plot(x,y)
> dev.off
```

Other devices you can use, with their corresponding file format are `pdf()`, `postscript()`, `png()` and `bitmap`.

Chapter 8

Lattice or Trellis Graphics

The lattice package in R provides an alternative to the base graphics system, it is an implementation of the ideas developed and implemented by Rick Becker and Bill Cleveland in the Trellis graphics system for the S language. Trellis displays were developed as a framework to make it easy the display of the relationship between a dependent variable and multiple factors. The lattice package allows you to do some wonderful things, however it's not always easy to use, and its syntax differs from the R base graphics system. Since lattice is mostly compatible with the Trellis graphics system in S-Plus, the following documents written for Trellis probably provide the best introduction to lattice and to the concept of trellis displays:

- S-PLUS Trellis Graphics User's Manual [2]
- A Tour of Trellis Graphics [4]

they can both be downloaded at the following URL:

<http://cm.bell-labs.com/cm/ms/departments/sia/project/trellis/software.writing.html>

As usual the definitive reference is provided by the (especially in this case not very beginner-friendly) package documentation:

- Lattice Reference Manual
<http://cran.r-project.org/doc/packages/lattice.pdf>

the following articles also contain useful info:

- Some notes on lattice [17]
<http://www.ci.tuwien.ac.at/Conferences/DSC-2003/Proceedings/Sarkar.pdf>
- R Lattice Graphics [14]
<http://www.ci.tuwien.ac.at/Conferences/DSC-2001/Proceedings/Murrell.pdf>
- Lattice, an implementation of Trellis graphics in R [16]
http://cran.r-project.org/doc/Rnews/Rnews_2002-2.pdf

- The Visual Design and Control of Trellis Display [3]
<http://cm.bell-labs.com/stat/doc/trellis.jcgs.col.ps>
- Trellis Display: Modelling Data from Designed Experiments [5]
<http://cm.bell-labs.com/stat/doc/dae.trellis.col.ps>

8.1 Overview of Lattice Graphics

Table 8.1: Lattice Graphics Functions

Function
xyplot
barchart
dotplot
stripplot
bwplot

8.2 Introduction to model formulae and multi-panel conditioning

We will use the ‘rats_trellis.txt’ dataset to illustrate how conditioning based on one or more factors work in lattice. The dataset contains data from a fictitious experiment in which a researcher is investigating the effects of alcohol and drug consumption on social interactions in rats. The researcher has studies two species of rats (Kalamani vs Yuppy), each rat has been observed in four experimental conditions, given by the combination of two factors: administration of drug (Drug vs No-Drug) and administration of alcohol (Al vs No-Al). The dependent variable is the number of social interactions observed in each of the four conditions. We’ll read in the data first:

```
> dats <- read.table("rats_trellis.txt", header=TRUE)
> dats$subj <- as.factor(dats$subj)
> head(dats)
  subj socialint alcohol  drug species
1    1         7      Al   Drug  Yuppy
2    1         6 No-Al   Drug  Yuppy
3    1         6      Al No-Drug Yuppy
4    1         4 No-Al No-Drug Yuppy
5    2         5      Al   Drug  Yuppy
6    2         4 No-Al   Drug  Yuppy
```


we'll start visualising the data along one dimension, the species. This is a single dimension, that can be easily handled by the base graphics system, but can be equally well displayed with lattice. Since the high level lattice plotting functions require the data to be entered as a dataframe, we'll use the `aggregate` function to get a dataframe with the mean values of `socialint`, the dependent variable, on the bases of the species:

```
> bySpec <- aggregate(dats$socialint,
+                     by=list(species=dats$species), FUN=mean)
> names(bySpec)[which(names(bySpec)=="x")] <- "socialint"
> bySpec
  species socialint
1 Kalamani    5.3125
2   Yuppy    5.6875
```

lattice uses a model formula syntax, you give it a dataframe, and specify how you want a variable to be displayed along the dimensions of one or more factors. In our case, we want `socialint ~ species` (you could read the `~` “as explained by species”), the barchart can be produced with the code below, and is displayed in Figure 8.1

```
> trellis.device()
> p11 <- barchart(socialint ~ species, data=bySpec)
> print(p11)
```

Now suppose we want to visualise the relationship between `socialint` and `species` on the bases of alcohol administration. We could achieve also this with the barplot in the base graphics system. In lattice there are two ways of doing it, one is by the use of a “grouping” factor, this yields a display similar to the barplot function. The second way is by multi-panel factor conditioning. The advantage of multi-panel conditioning, as we will see soon, is that it can be extended to an unlimited number of factors. Let's start with the first solution, in which we use a grouping factor. First we create a suitable dataframe:

```
> bySpecA1 <- aggregate(dats$socialint,
+                       by=list(species=dats$species,
+                               alcohol=dats$alcohol), FUN=mean)
> names(bySpecA1)[which(names(bySpecA1)=="x")] <- "socialint"
```

then we produce the barchart (Figure 8.2

```
> p12 <- barchart(socialint ~ species, groups=alcohol,
+                data=bySpecA1, auto.key=TRUE)
> print(p12)
```

Figure 8.1: Social Interactions by Species

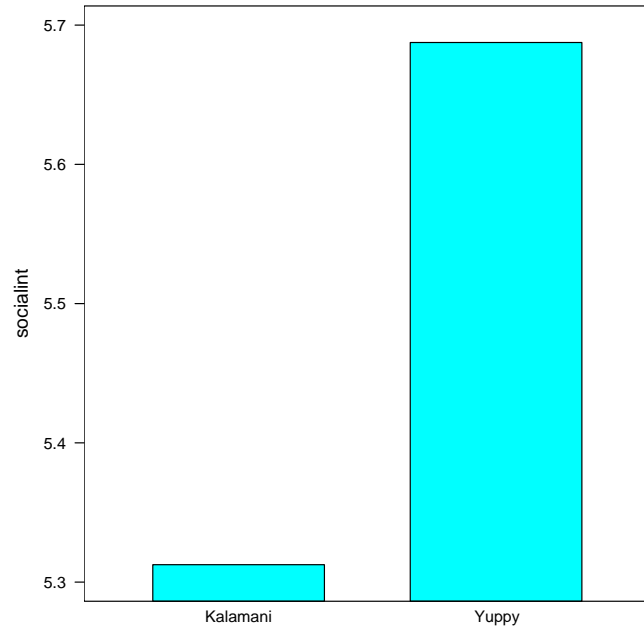
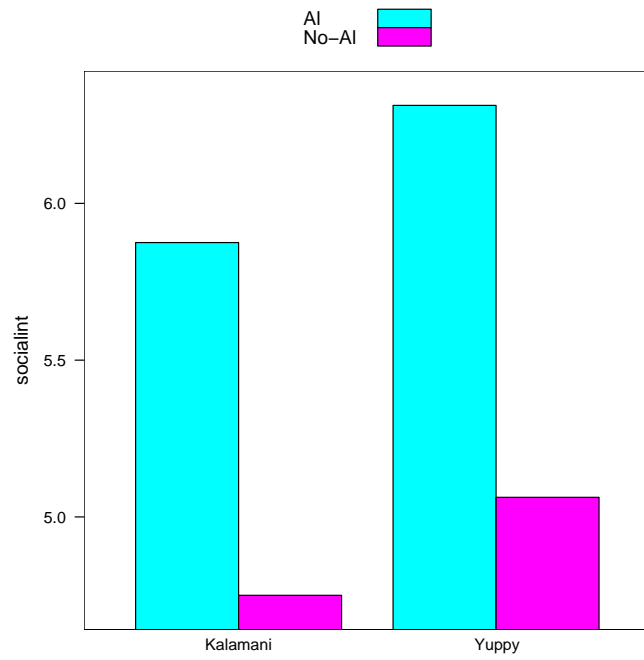


Figure 8.2: Social Interactions by Species and Alcohol Administration with Grouping Factor

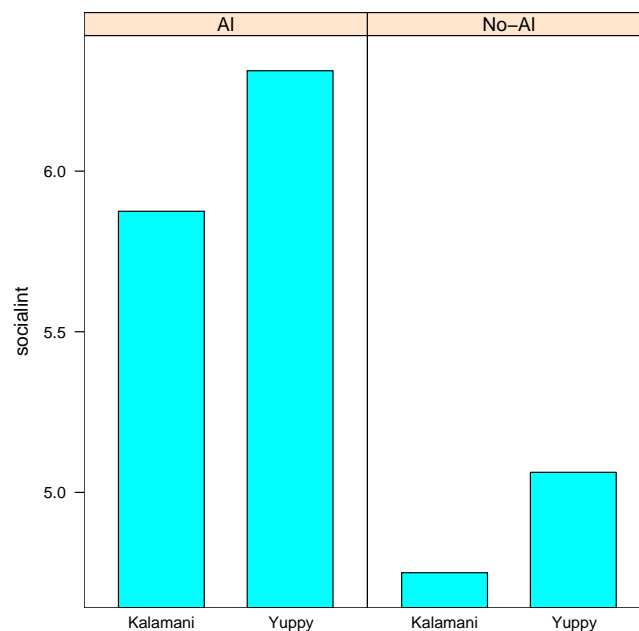


the grouping factor is given by the `groups` argument, notice also that we have set `auto.key` to `TRUE` in order to add an automatic legend. the second way of doing the graph is by multi-panel conditioning, we achieve this by putting `alcohol` as a conditioning factor with the `|` syntax.

```
> p13 <- barchart(socialint ~ species | alcohol, data=bySpecAl)
> print(p13)
```

the resulting graph is displayed in Figure 8.3 we could have swapped `species`

Figure 8.3: Social Interactions by Species and Alcohol Administration with Multi-panel Conditioning



for `alcohol` as the conditioning factor, which way gives the more effective display depends from case to case, and it's up to the user to decide. Finally we'll consider the case in which all factors are included in the display, this cannot be easily achieved with traditional graphics, but it is easily done in lattice, we just add `drug` to the conditioning factors (Figure 8.4

```
> bySpecAlDrug <- aggregate(dats$socialint,
+                           by=list(species=dats$species,
+                                   alcohol=dats$alcohol, drug=dats$drug),
+                           FUN=mean)
> names(bySpecAlDrug)[which(names(bySpecAlDrug)=="x")] <- "socialint"
> p14 <- barchart(socialint ~ species | alcohol * drug,
+                data=bySpecAlDrug)
> print(p14)
```

we could have used a grouping variable also in this case rather than using two conditioning factors (Figure 8.5)

```
> p15 <- barchart(socialint ~ species | alcohol,
+                 groups=drug, data=bySpecAlDrug,
+                 auto.key=TRUE)
> print(p15)
```

again what is the best display is up to the user to decide and depends from case to case.

8.3 barchart

We will look at an example of a barchart display of a dependent measure conditional on three factors. The dataset ‘line_matching.txt’ contains data on an imaginary experiment in which a psychophysician wants to measure the accuracy of matching the length of a segment for three groups of people (Gr. 1, Gr. 2, Gr. 3) for segments of four different lengths (L1, L2, L3, L4). The third factor the psychophysician is interested in is whether matching accuracy changes depending on the colour (blue vs red) of the segment to be matched, he measures this as a within subjects factor. The matching accuracy is measured as the error, or displacement (positive or negative) from the actual segment length. The dataset contains the mean values for the three groups. Below is the code for producing the barchart, the resulting plot can be seen in Figure 8.6:

```
> dats <- read.table("line_matching.txt", header=TRUE)
> trellis.device()
> oldpar <- trellis.par.get("superpose.polygon")
> trellis.par.set(superpose.polygon =
+ list(col = c("darkslateblue", "indianred")))
> myGraph <- barchart(error ~ length | group, groups=color,
+ data=dats, origin=0, ylab="Error (cm)",
+ xlab="Segment Length",
+ auto.key=TRUE, as.table=TRUE)
> print(myGraph)
> trellis.par.set(superpose.polygon = oldpar)
```

We’re showing the bars for the two levels of the “colour” factor side by side in the same panel, this is done by using the factor in the `groups` argument. A different display could have been achieved by putting the “colour” factor as an additional conditioning variable

```
> barchart(error ~ length | group * color,
+ data=dats, origin=0, ylab="Error (cm)", xlab="Segment Length",
+ auto.key=TRUE, as.table=TRUE)
```

Figure 8.4: Social Interactions by Species, Alcohol Administration and Drug with Multi-panel Conditioning

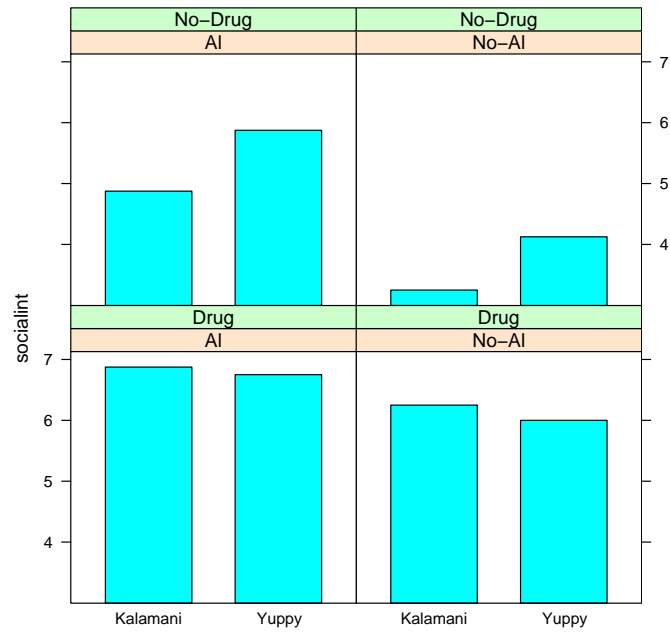
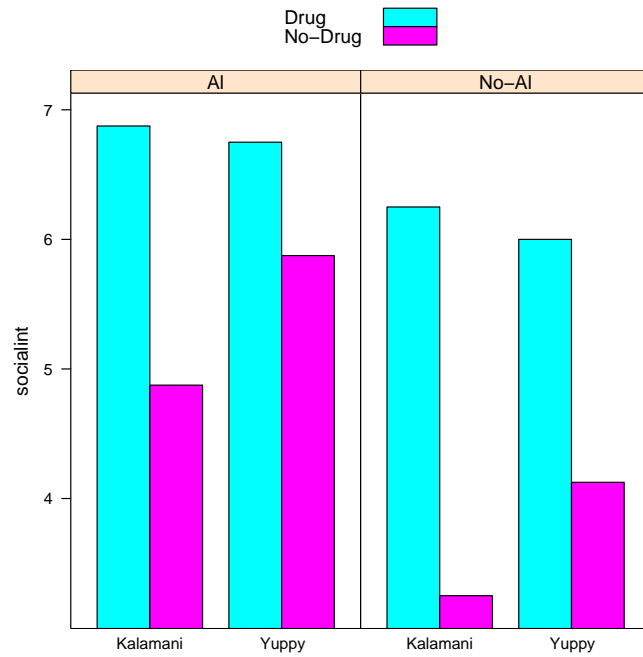


Figure 8.5: Social Interactions by Species, Alcohol Administration and Drug with Grouping Factor



in this case the bars for each level of the factor would have been drawn in different panels (the number of panels would have doubled).

The fill colour for the bars can be modified changing the colour option for `superpose.polygon`.

The great thing about trellis graphics is that they allow you to display the relationship between a dependent variable and multiple factors seamlessly. Suppose that, continuing the above example, the psychophysician has tested the line matching accuracy on the three groups both before and after a period of visuo-motor training. The dataset containing the data with this new factor is in the file 'line_matching_training.txt'. We just need to add the new factor `session` (pre-training vs post-training) to the conditioning variables to obtain the new plot. The modified call to the `barchart` function is shown below and the resulting graph can be seen in Figure 8.7.

```
> myGraph <- barchart(error ~ length | group * session, groups=color,
+                   data=dats, origin=0, ylab="Error (cm)",
+                   xlab="Segment Length",
+                   auto.key=TRUE, as.table=TRUE)
> print(myGraph)
```

8.4 dotplot

8.5 histogram

```
> histogram(dats$adj)
> histogram(~dats$adj|dats$congr)
```

8.6 Interaction Plots

Example of interaction plot with 3 factors:

```
> print(bwplot(dats$lat~dats$congr|dats$acc,groups=dats$isi,
+             panel='panel.superpose',panel.groups='panel.linejoin',
+             auto.key=list(points=FALSE,lines=TRUE,space='top'),
+             scales=list(cex=.8),ylim=c(280,400),
+             ylab='Medie delle latenze in ms',xlab='SOA'))
```

8.7 Customising Lattice Graphics

8.7.1 Textual Elements

Strip Labels

The labels of the panels strips are the names of the levels of the conditioning factor variable. To change the labels just change the names for the factor levels:

```
> levels(dats$acc)
[1] "e" "h"
> levels(dats$acc)[1]
[1] "e"
> levels(dats$acc)[1]<-'imprecise'
> levels(dats$acc)[2]<-'corrette'
> levels(dats$acc)
[1] "imprecise" "corrette"
```

8.8 Writing Panel Functions

8.8.1 Combining Panel Functions

Rather than writing a new panel function from scratch, often you just want to add some elements to a plot, for example a regression line, or error bars. The easiest and probably best way to do this is writing a panel function that combines two standard panel functions. There is a number of predefined panel functions (see `?panel.functions`) that can be used to add lines, grids etc. . . , to a scatterplot, barchart or other higher level plotting lattice function.

We will start with a very simple example, adding a horizontal line at a fixed height in a dotplot. The 'line_matching' dataset described in the barchart example, can be very well visualised also through a dotplot (Figure 8.8)

```
> oldpar <- trellis.par.get("superpose.symbol")
> trellis.par.set(superpose.symbol =
+                 list(col = c("darkslateblue",
+                               "indianred"), pch=19))
> dotplot(error ~ length | group, groups=color,
+         data=dats, origin=0, ylab="Error (cm)",
+         xlab="Segment Length", auto.key=TRUE,
+         aspect=1, as.table=TRUE)
```

however, since the data represent positive or negative displacements from zero, it would be nice to add a horizontal line passing from zero. In order to

have this, we will write a panel function that combines the `panel.dotplot` function with the `panel.abline` function that we'll use to add the horizontal line

```
> panel.hRefDotplot <- function(x, y, ref=NULL, ...){
+   panel.dotplot(x, y, ...)
+   panel.abline(h=ref, ...)
+ }
```

our new `panel.hRefDotplot` panel function accepts three arguments, `x` and `y`, which are the “standard” arguments given by the higher level plotting functions like `dotplot` to panel functions to specify the data to draw. The third argument represents the position at which to draw the horizontal line of reference for the data, we want it to be zero in this case, but passing the argument as a variable rather than hard-coding the value into the panel function will allow us to recycle this panel function in case we want the horizontal reference line drawn at some other points in the future. Besides these arguments, our panel function accepts also an undefined number of other arguments, which are designated by the ‘...’ notation. These are usually graphics parameters that can be specified in the high level plotting function. The contents of our `panel.hRefDotplot` function are very simple, we call first `panel.dotplot` to draw the standard dotplot, and then we call `panel.abline` giving it the value of `ref` to draw the horizontal line in each panel. The actual plot is done by calling the high level `dotplot` function specifying `panel.hRefDotplot` as the panel function to use

```
> dotplot(error ~ length | group, groups=color,
+ data=dats, origin=0, ylab="Error (cm)",
+ xlab="Segment Length", auto.key=TRUE, aspect=1,
+ as.table=TRUE, ref=0, panel=panel.hRefDotplot)
```

notice the last line of the call, first we're telling `dotplot` to use our `hRefDotplot` function to do the plotting by specifying the `panel` argument, second we're specifying another argument, `ref` in the call, this is not a standard argument, but it will be automatically passed to our panel function to decide at which height to draw the horizontal line. The resulting plot can be visualised in Figure 8.9.

Figure 8.6: Line Matching Barchart

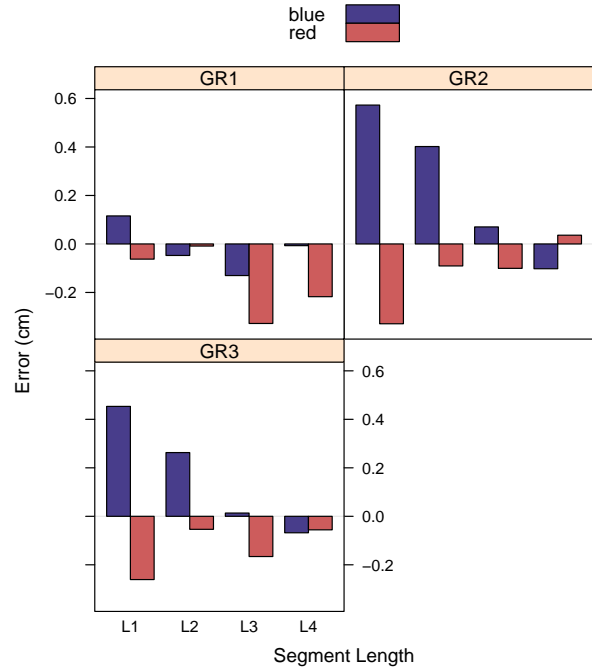


Figure 8.7: Line Matching Training Barchart

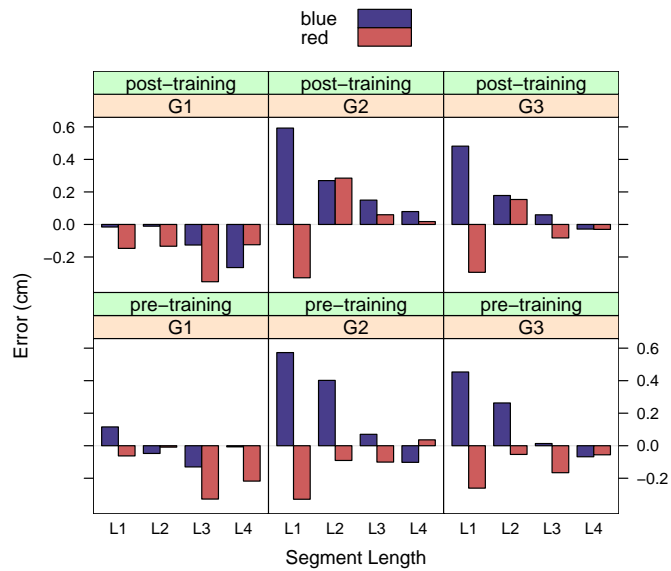


Figure 8.8: Dotplot of the Line Matching Dataset

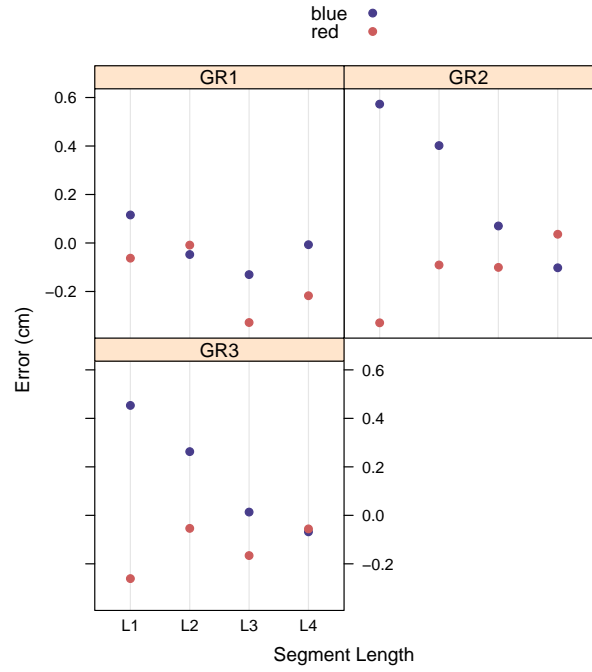
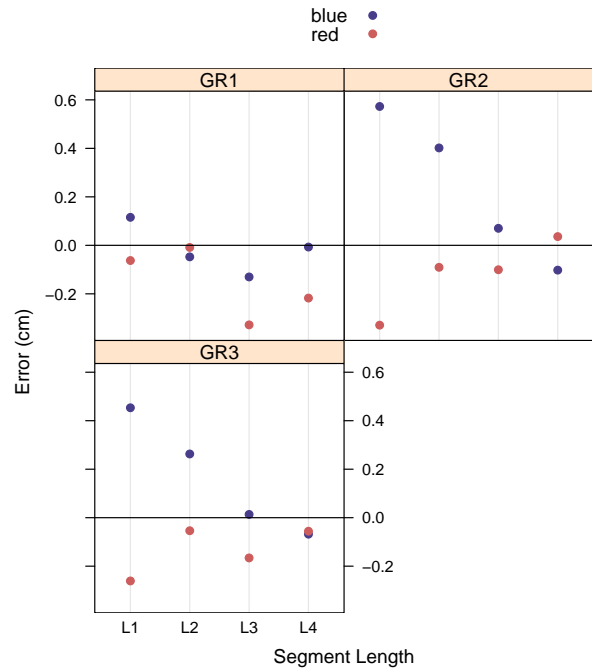


Figure 8.9: Dotplot of the Line Matching Dataset, with Horizontal Reference Line



Chapter 9

Probability Distribution Functions

9.1 The Bernoulli distribution

A random variable X that takes a value of 0 or 1 depending on the result of an experiment that can have only two possible outcomes, follows a Bernoulli distribution. If the probability of one outcome is p , the probability of the other outcome will be $p - 1$:

$$P(X = 1) = p$$

$$P(X = 0) = p - 1$$

The expected value of a Bernoulli random variable is p :

$$E[X] = 1 \cdot p + 0 \cdot p = p$$

the variance of a Bernoulli random variable is given by:

$$Var(X) = E[X^2] - E[X]^2 = 1^2 \cdot p + 0^2 \cdot p - p^2 =$$

$$Var(X) = p - p^2 = p \cdot (1 - p)$$

As far as I know, there are no special functions in R for computations related to the Bernoulli distribution (density, distribution, and quantile function). However, since the Bernoulli distribution is a special case of the binomial distribution, with parameter *size* equal to 1, the R functions for the binomial distribution can be used for the Bernoulli distribution as well. Most of the calculations involved are quite simple anyway. The density function can be calculated as follows:

```
> dbinom(x=0, size=1, prob=0.7)
[1] 0.3
> dbinom(x=1, size=1, prob=0.7)
[1] 0.7
```

The cumulative distribution function can be computed as follows:

```
> pbinom(q=0, size=1, prob=0.7)
[1] 0.3
> pbinom(q=1, size=1, prob=0.7)
[1] 1
```

The quantile function can be computed as follows:

```
> qbinom(p=0.3, size=1, prob=0.7)
[1] 0
> qbinom(p=1, size=1, prob=0.7)
[1] 1
```

Finally, one can generate a random sample from a Bernoulli distribution as follows:

```
> rbinom(n=10, size=1, prob=0.7)
[1] 1 0 1 0 1 1 0 1 0 1
```

9.2 The binomial distribution

9.3 The normal distribution

pnorm

You can use the `pnorm` command to find out the probability value associated with a given z point in the normal distribution. For example:

```
> pnorm(1.96)
```

will give the value of the area under the normal distribution curve from $-\infty$ to 1.96.

qnorm

The `qnorm` command is the inverse of the `pnorm` command, in that it gives the z point associated with a given probability area under the normal curve. For example:

```
> qnorm(0.975)
```

dnorm

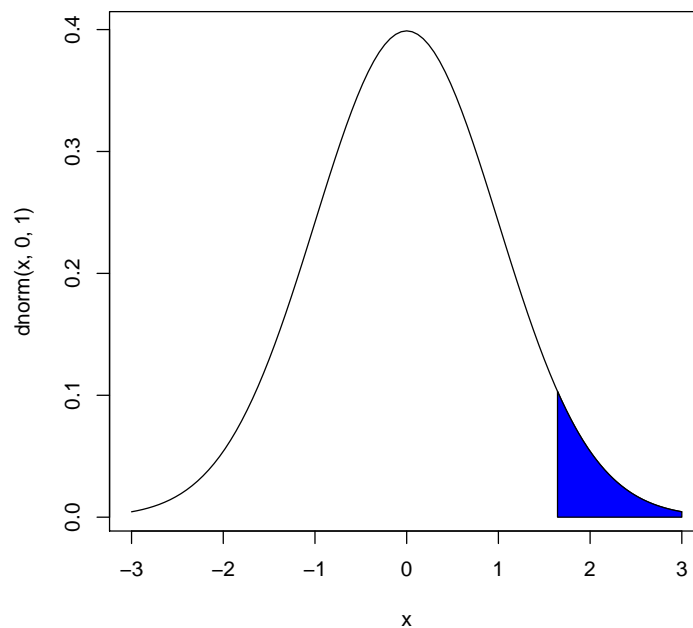
The `dnorm` function provides the density function for the normal distribution. Using the `dnorm` function we can for example make a nice plot of a normal distribution with mean equal to 0 and a standard deviation of 1 and colour its extreme right 0.05% tail with the following code. The result is displayed in Figure 9.1

```

> curve(dnorm(x,0,1), from=-3, to=3)
> coord <- seq(from=0+qnorm(.95)*1, to=3, length=30)
> dcoord <- dnorm(coord, 0, 1)
> polygon(x=c(0+qnorm(.95)*1, coord, 3),
+         y=c(0, dcoord, 0), col = "blue")

```

Figure 9.1: The Normal Distribution



We can use the `pnorm` function to test hypotheses with a z-test. Suppose we have a sample of 40 computer science students with a mean short term memory span of 7.4 digits (that is, they can repeat in sequence, about 7 digits you read them, without making a mistake), and standard deviation of 2.3. The mean for the general population is 6.5 digits, and the variance in the population is unknown. We're interested in seeing if the memory span for computer science students is higher than that of the general population. We'll run a z-test as

$$z = \frac{\bar{x} - \mu_{\bar{x}}}{\frac{s}{\sqrt{n-1}}}$$

where \bar{x} is the mean short term memory span for our sample, and s is its standard deviation.

```

> z <- (7.4-6.5)/(2.3/sqrt(40-1))
> z

```

```

[1] 2.443695
> pnorm(z)
[1] 0.9927311
> 1-pnorm(z)
[1] 0.007268858 #this would do for a one tailed test
> (1-pnorm(z))*2
[1] 0.01453772

```

so the z value for our sample is 2.443695, we then get the area under the curve from $-\infty$ to our z-value, which gives us the probability of a score lower than 7.4. To get the p-value we subtract that probability value from 1, this would give us the probability of getting a score equal to or higher than 7.4 for a one tailed test. Since we want a two-tailed test instead, we multiply that value by two, to get our p-value, which is nonetheless significant.

The `qnorm` command on the other hand, can be used to set confidence limits on the mean, for the example above we would use the following formula:

$$CI = \bar{x} \pm z_{\alpha/2} \frac{s}{\sqrt{n-1}}$$

to set a 95% confidence interval on the mean short term memory span for the computer science students:

```

> alpha <- 0.05
> s <- (2.3/sqrt(40-1))
> zp <- (1-alpha/2)
> ciup <- 7.4+zp*s
> cilow <- 7.4-zp*s
> cat("The 95% CI is: \n",cilow,"<","mu","<","ciup","\n")
The 95% CI is:
7.040913 < mu < 7.759087

```

The following code summarises the situation graphically, as shown in Figure 9.2

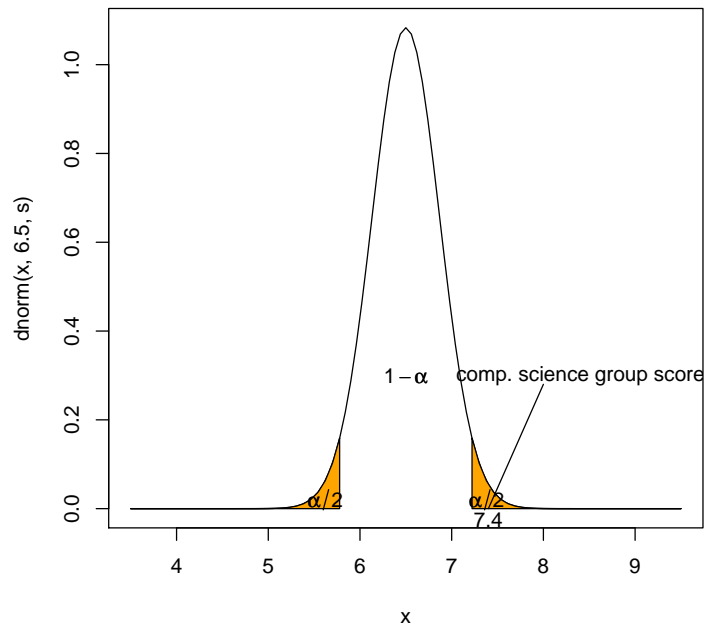
```

> s <- (2.3/sqrt(40-1))
> up <- seq(from = 6.5+qnorm(.975)*s, to = 9.5, length = 30)
> low <- seq(from =3.5 , to =6.5-qnorm(.975)*s , length = 30)
> dup <- dnorm(up,6.5,s)
> dlow <-dnorm(low,6.5,s)
> curve(dnorm(x,6.5,s), from= 3.5, to= 9.5)
> polygon(x = c( 6.5+qnorm(.975)*s, up, 9.5),
+         y = c(0, dup, 0), col = "orange")
> polygon(x = c(3.5,low, 6.5-qnorm(.975)*s),
+         y = c(0, dlow, 0), col = "orange")
> text(x=c(7.38,5.62),y=c(0.02,0.02),expression(alpha/2))

```

```
> text(x=c(6.5),y=c(0.3),expression(1-alpha))
> lines(x=c(8,7.4),y=c(0.28,0))
> text(x=c(7.4),y=c(-0.025),expression(7.4))
> text(x=c(8.4),y=c(0.3),"comp. science group score")
```

Figure 9.2: Short Term Memory Span Experiment, z-test



Chapter 10

Hypothesis Testing

10.1 χ^2 test

Testing Hypotheses about the Distribution of a Categorical Variable

You can use the χ^2 test to verify hypotheses about the distribution of a categorical variable, for example you might want to know whether, as far as handedness is concerned, your experimental sample is representative of the general population and follows a distribution of 80% right-handed, 15% left-handed, and 5% using both hands indifferently¹. The frequencies for your sample are given in Table 10.1 in R you can use the `chisq.test()`

Table 10.1: Handedness Frequencies for an Experimental Sample

Right	Left	Both	Tot.
93	18	2	113

function to see if the distribution of your sample differs from that of the general population:

```
> chisq.test(x=c(93,18,2), p=c(0.8,0.15,0.05))
```

```
Chi-squared test for given probabilities
```

```
data: c(93, 18, 2)
X-squared = 2.4978, df = 2, p-value = 0.2868
```

you give the function a vector `x` with the frequencies for each category in your sample, and a vector `p` of the theoretical probability for each category, R returns the value of the χ^2 statistics and its associated p-value. In our case the p-value indicates that the χ^2 statistics is not significant, assuming

¹Fake data

$\alpha = 0.05$, so we can say that the distribution for the handedness variable in our sample does not differ from that of the general population.

Testing Hypotheses about the Association between Two Categorical Variables

The χ^2 test can also be used to verify a possible association between two categorical variables, for example sex and cosmetic products usage (classified as “high”, “medium” or “low”). In this case the data are best summarised by a contingency table as Table 10.2 which presents the data on cosmetic usage for a sample of 44 males and 67 females. The formula for the χ^2

Table 10.2: Cosmetic Usage by Sex

	High	Medium	Low	Tot.
Males	6	11	27	44
Females	25	30	12	67
Tot.	31	41	39	111

statistics is

$$\chi^2 = \sum \frac{(f_e - f_o)^2}{f_e} \quad (10.1)$$

where f_e are the expected frequencies and f_o are the observed frequencies, to get the statistics and perform a test of significance on it with R you can first arrange the data in a matrix:

```
> cosm <- c(6,11,27,25,30,12)
> cosm <- matrix(cosm, nrow=2, byrow=TRUE)
```

then you can directly give the matrix to the `chisq.test()` function to test the null hypothesis that there is no association between sex and cosmetic usage

```
> chisq.test(cosm)
```

Pearson's Chi-squared test

```
data:  cosm
X-squared = 22.4159, df = 2, p-value = 1.357e-05
```

the p-value tells us that the null hypothesis does not hold, there is indeed an association between sex and degree of cosmetic usage. If you store the results of the test in a variable, you'll be able to get the matrices of the expected frequencies and residuals

```
> cosmtest <- chisq.test(cosm)
```

```

> cosmtest$expect ## expected frequencies
>
      [,1] [,2] [,3]
[1,] 12.28829 16.25225 15.45946
[2,] 18.71171 24.74775 23.54054
> cosmtest$resid^2
      [,1] [,2] [,3]
[1,] 3.217907 1.697374 8.615054
[2,] 2.113252 1.114694 5.657647
> sum(cosmtest$resid^2)
[1] 22.41593

```

If you have only a few observations and want to use the Yate's correction for continuity you have to set the optional argument `correct` to true

```
> chisq.test(x=mydata, correct=TRUE)
```

Very often you don't already have your data tabulated in a contingency table, however if you have a series of observations classified according to two categorical variables you can easily get a contingency table with the `table()` function. The next example uses the data in the file `'hair_eyes.txt'`, that lists the hair and eyes colours for 118 people. We will first build our contingency table from that:

```

> dats <- read.table("hair_eyes.txt", header=TRUE)
> tabdats <- table(dats$hair, dats$eyes)
> tabdats
      eyes
hair  blue brown green
blonde  16    9    6
brown   14   32    9
red     9    7   16

```

now we can perform the χ^2 test directly on the table to see if there is an association between hair and eyes colour

```
> chisq.test(tabdats)
```

Pearson's Chi-squared test

```

data:  tabdats
X-squared = 21.8622, df = 4, p-value = 0.0002135

```

the χ^2 statistics is highly significant, so there clearly is an association between hair and eyes colour.

10.2 Student's t test

10.2.1 One sample t-test

To test if a group's mean is different from a given population's mean (μ), you can run a simple t-test in R. The syntax is as follows:

```
> t.test(gr1, mu = 17)
```

where `gr1` is the vector containing the data of the group, and `mu` is an optional argument ² specifying the population's mean.

- `mu = 0` [Default] or another value
mu is the value specifying the population's mean.
- `alternative = "two.sided"` [default] or "less" or "greater"

this option allows choosing whether the test should be two-tailed or one-tailed. In particular the alternative hypothesis H_1 becomes:

- "two.sided", true mean of the group is not equal to mu
- "less", true mean of the group is lesser than mu
- "greater", true mean of the group is greater than mu

- `conf.level = 0.95` [Default] or another value
this option sets the confidence level for the estimation of the confidence interval.

10.2.2 Two samples t-test

With R you can easily run a t test to find out whether the difference between the means of two groups is significant. The syntax is as follows:

```
> t.test(gr1, gr2)
```

where `gr1` is the data vector of the first group and `gr2` is the data vector of the second group. The following options can also be selected:

The following options can also be selected:

- `mu = 0` [Default] or another value
mu is the value given by the difference between the means under H_0 , in other words it is the value against which you test the alternative hypothesis.

²If you leave it out it defaults to 0

- `alternative = "two.sided"` [default] or `"less"` or `"greater"`

this option allows choosing whether the test should be two-tailed or one-tailed. In particular the alternative hypothesis H_1 becomes:

- `"two.sided"`, true mean of the group is not equal to mu
- `"less"`, true mean of the group is lesser than mu
- `"greater"`, true mean of the group is greater than mu

- `paired = FALSE` [Default] or `TRUE`

If the `TRUE` option is selected, then a *paired t-test* is performed, note that in this case the number of observations has to be equal in the two groups. Missing values are removed if `paired` is `TRUE`.

A shortcut to choose a paired t-test is:

```
> t.test(gr1-gr2)
```

- `var.equal = FALSE` [Default] or `TRUE`

Tells whether the variances of the two groups should be treated as equal. If `TRUE` is selected, then the pooled variance is used to calculate the variance. If the option is left to default or `FALSE` is selected, then the Welch (or Satterthwaite) approximation to the degrees of freedom is used.

- `conf.level = 0.95` [Default] or another value

this option sets the confidence level for the estimation of the confidence interval.

10.3 The Levene Test for Homogeneity of Variances

The Levene test can be used to test if the variances of two samples are equal or not. This procedure is important because some statistical tests, such as the two independent sample t test, are based on the assumption that the variances of the samples being tested are equal. In R the function to perform the Levene test is contained in the `car` package, so in order to call it the package must be installed and the `car` library has to be loaded with the command:

```
> library(car)
```

The syntax to run a Levene test is:

```
> levene.test(y, group)
```

where `y` is a vector with the values of the dependent variable we are measuring, and `group` is another vector defining the group to which a given observation belongs to. With an example this will be clearer. Let's imagine that in a verbal memory task we have the measures (number of words recalled) of three groups that have been given three different treatments (three different rehearsal procedures). We want to see if the variances of the three groups are equal. The data are stored in a text file "`verbal.txt`", as in the table 10.3, but with no headers whatsoever.

Table 10.3: Data for the Levene test example

Procedure 1	Procedure 2	Procedure 3
12	13	15
9	12	17
9	11	15
8	7	13
7	14	16
12	10	17
8	10	12
7	10	16
10	12	14
7	9	15

First we will create the vector with the values of the dependent variables, by reading in the file with the `scan` function:

```
> values <- scan("verbal.txt")
```

Now we need to create the second vector defining the group to which a given observation on the data vector belongs to. We will use the `rep` function to do this, and for convenience will give the labels "1", "2" and "3" to the three groups:

```
> groups <- as.factor (rep(1:3, 10))
```

Notice that we have defined this vector as a factor, this was necessary because we were using numerical labels, it wouldn't have been if we had used a character or a string labels. Now the data are in a proper format for using the `levene.test` function. This format can be called "one row per observation", and we will encounter it often later when we talk about ANOVA, you can find more examples and clarifications on this format there. To perform the test we can use the following commands:

```
> levene.test(values, groups)
```

The R output is given below, as you can see we have an F value with its associated p value, if the F value is significant, then the variances in the groups are not equal. In our case below the p value is far greater than 0.05 so we can safely accept the null hypothesis that the variances in the three groups are equal.

Levene's Test for Homogeneity of Variance

	Df	F value	Pr(>F)
group	2	0.3391	0.7154
	27		

Chapter 11

Correlation and Regression

Let's imagine we want to see if there is a correlation between a manual and an oculomotor reaction time (RT) task. The manual RT task requires to press one of two buttons depending on the colour assumed by the fixation point after a variable delay. The oculomotor RT task requires to make a saccade towards the right or the left depending on the colour taken by the fixation point. The mean RTs in milliseconds, for 37 subjects are in the file `m_o_rt.txt`. We read in the data, and then calculate the correlation for the sample:

```
> dats <- read.table("m_o_rt.txt", header=TRUE)
> cor(dats$man,dats$ocul)
[1] 0.5030888
> cor(dats$man,dats$ocul)^2 ## R squared
[1] 0.2530983
```

the correlation between the two measures in the sample is moderate, and it would account for about 25% of the variance. We need also to test if this correlation could be extended to the population

```
> cor.test(dats$man,dats$ocul)
```

Pearson's product-moment correlation

```
data: dats$man and dats$ocul
t = 3.4439, df = 35, p-value = 0.001505
alternative hypothesis: true correlation is not equal to 0
95 percent confidence interval:
 0.2139450 0.7111784
sample estimates:
      cor
0.5030888
```

the p-value tells us that we can extend our finding in the population, so people with a quick hand would also have a quick eye.

11.1 Linear Regression

Regression models are expressed in R with a symbolic syntax that is clean and convenient. If you have a dependent variable `rts` and you want to check the influence of an independent variable `age` on it, this would be expressed as

```
> rtfit <- lm(rts~age)
> plot(rts~age)
> abline(rtfit)
```

if you want to add another predictor `score`

```
> rtfit2 <- lm(rts~age+score)
```


Chapter 12

ANOVA

12.1 One-Way ANOVA

We'll start directly with an example to show how to perform a one-way ANOVA. Let's imagine an experimenter is studying the effects of sleep deprivation on verbal memory. He selects a random sample of healthy subjects, and assigns them to three treatments: 4 hours, 6 hours and 8 hours of total sleep. The subjects are then given a memory test, consisting of a long list of words to remember after a minute of delay from their presentation. The dependent variable is the number of words correctly recalled in a minute time. The data are shown in table 12.1:

Table 12.1: Data for the one-way ANOVA example

4 hours	6 hours	8 hours
12	13	15
9	12	17
9	11	15
8	7	13
7	14	16
12	10	17
8	10	12
7	10	16
10	12	14
7	9	15

The data are stored in a text file "sleepdep.txt" in the above format, but without any header. First we read in the data with the `scan` function and assign them to a vector that we'll call `recalled`:

```
recalled <- scan("sleepdep.txt")
```

then we need to create another vector which holds the *level* of the factor “hours of sleep” for each observation, this vector can contain either numbers (to be assigned as factors), or single characters such as “c” or “d”, but not strings such as “cod1”, “cod2” and so on. Since our data, stored in the vector “recalled” are organised in an ordered sequence, with observations belonging to treatment “4 h, 6 h, 8 h ...” repeated 10 times, we need to mimic this structure to create our new vector. We’ll use the `rep` function to do this, we’ll call our levels “1, 2 and 3” for sake of simplicity.

```
> treatment <- as.factor (rep(1:3), 10))
```

now we’ll put the two vectors together into a data frame:

```
> sleepdata <- data.frame (recalled, treatment)
```

well, the data are now in a proper format to perform our analysis of variance, we can call this format “**one row per observation**”, since in each row of our data frame we hold the value of the dependent variable in one column and the level of the factor that we are manipulating in the other column for a single observation. As for the analysis, we go like this:

```
> oneway.test (sleep~treat, data=sleepdata, var.equal=TRUE)
```

the `summary` command tells R to print out a nicely formatted summary with the results of our analysis. The actual command that performs the analysis is `aov`, and what follows this command is the specification of our model for the analysis. In this case we want to see if the number of words recalled depends on the treatment, the `(recalled ~ treatment)` statement does just this, because the tilde “~” means “explain recalled on the basis of treatment”. Finally, the `data` statement specifies the object in which are stored the data for this analysis.

Below there is the summary produced by R for this analysis.

```
              Df  Sum Sq  Mean Sq F value Pr(>F)
treatment     2   194.867   97.433  27.838 2.746e-07 ***
Residuals    27    94.500    3.500
---
```

```
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

We have the sums of squares for the treatment effect and for the error term and the F value for the treatment effect with its significance value. Since the p value in this case is very, very small, it is written in scientific notation. However we can use the significance codes given at the bottom of the print out to interpret the p value as very close to zero, at least smaller than 0.001, so it is highly significant.

12.1.1 One-Way ANOVA with Unequal Sample Sizes

The commands to perform an univariate ANOVA with unequal sample sizes in R are the same as in the case the samples have the same size. If the data are not in the “one row per observation” format however, the details of the procedures to use to format them properly are slightly different. Let’s see an example. We want to see if rats run faster in galleries painted in red, rather than galleries painted in green or blue. We wanted to run 20 subjects for each condition in this experiment, however for various reasons 8 out of the total number of rats, 3 from the “red” group, 4 from the green group and “1” from the blue group, had to be dropped from further analyses. The data are shown in table 12.2, the time taken to run all over the gallery for each subject is shown in seconds. The data are store in a file called “ratsrun.txt”.

Table 12.2: Data for the one-way ANOVA example with unequal sample sizes

Red	Green	Blue
120	133	115
90	127	163
130	113	121
103	126	104
145	97	132
127	105	144
138	107	121
165	102	104
104	126	142
123	143	123
120	135	124
142	126	143
130	111	122
103	132	134
115	141	96
122	107	105
128		126
		139
		112

We read in the data as usual:

```
> runtimes <- scan("ratsrun.txt")
```

Now we need to create a vector with the levels of the “colour” factor for each observation. We will call the three levels “r”, “g” and “b”.

```
> colour <- as.factor( c(rep(c("r","g","b"),16),"r","b","b","b"))
```

Now we can put all together in a data frame:

```
> data <- data.frame(runtimes,colour)
```

Now the data are in a proper format to perform the analysis. Well, this was a bit an extreme example, of course it would have been easier if the data in the text file would have been in another format. It was just to show you that once you master some basic functions like `rep` and `c` (concatenate) in R, you can laid out your data or vectors you need, pretty much as you like, and you can adapt even to the most adverse situations.

Since we have unequal sample sizes, we should pay particular attention that the variances of the three groups are homogeneous. To check for this we can use the Levene test for homogeneity of variances:

```
> library(car)
> levene.test(runtimes,colour)
Levene's Test for Homogeneity of Variance
      Df F value Pr(>F)
group  2  0.0562 0.9454
      49
```

as you can see from the print out the resulting F value is far from being significant, so we can safely accept the null hypothesis that the variances in the three groups are homogeneous. Now we can run the test:

```
> oneway.test(runtimes ~ colour, data=data, var.equal=TRUE)
```

One-way analysis of means

data: runtimes and colour

F = 0.2731, num df = 2, denom df = 49, p-value = 0.7622

The F statistics is not significant, so there is no difference in running times for rats with the three different colours.

12.2 Repeated Measures ANOVA

12.2.1 One Within Subjects Factor

The syntax for performing a repeated measures ANOVA is a little more complex than the syntax for fully randomised designs. In a repeated measures design we take into account the effects of the subjects on our measures, that is the fact that different subjects will have different different baseline

means on a given measure (e. g. on a reaction time test). In this way we are able to tell apart the variability given by inter individual differences between our subjects, from the variability due to the manipulation of one or more independent variables, with a view to identify the latter with more precision. A consequence of this procedure is that while with other designs we used a common error term, in a repeated measures design we have to use different error terms to test the effects we are interested in and we have to specify this in the formula we use with R for `aov`. A good explanation of repeated measures designs in R is given by Baron and Li (2003) [1], and what follows in this discussion is mainly inspired by their work.

We will start with a simple example of a repeated measures design with one within subject factor at three levels. Let's imagine we want to test the effects of three different colours on a simple detection task. The data are presented in table 12.3, and represent subjects' reaction times (measured with a button press) for the detection of squares of three different colours ("blue", "black" and "red"). Each subject was tested under all the three conditions. The data are store in the file "`rts.txt`", with the same format as that shown in the table, but without any header and without the column specifying the subject's number. We will first read in the data and apply

Table 12.3: Example for repeated measures ANOVA with one within subjects factor

Subj	Blue	Black	Red
1	0.120	0.132	0.102
2	0.096	0.103	0.087
3	0.113	0.134	0.109
4	0.132	0.147	0.123
5	0.124	0.139	0.124
6	0.105	0.115	0.102
7	0.109	0.129	0.097
8	0.143	0.150	0.119
9	0.127	0.145	0.113
10	0.098	0.117	0.092
11	0.115	0.126	0.098
12	0.117	0.132	0.103

the usual transformations to get the one row per observation format:

```
> data <- scan("rts.txt") #read in the data
> colour <- as.factor(rep(c("blue","black","red"),12))
> #create a factor for colours with 3 lev.
> subj <- as.factor(rep(1:12,each=3))
> #create factor for subjects
```

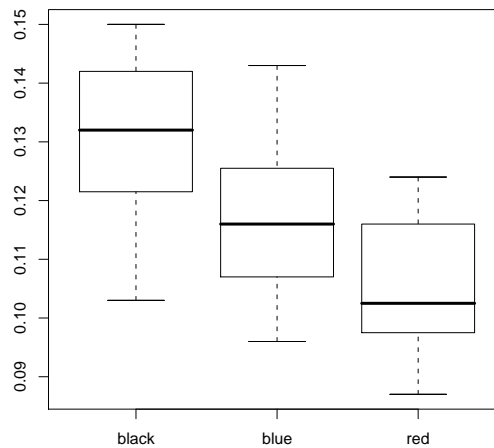
```
> table <- data.frame(data,colour,subj)
> #put everything in a dataframe
```

Now the data are ready for further analyses. We can first have a look at variability for the three colours by drawing a boxplot.

```
> boxplot(table$data ~ table$colour)
```

As you can see from fig. 12.1, the variability for the three conditions is pretty much the same and the three distributions seems to be approximately normal. The medians for the three distributions seems to differ, our analysis will tell us if these differences are significant.

Figure 12.1: Boxplot showing the distribution of RTs for the three colours



Below is shown the syntax for the analysis and with its output:

```
> summary(aov(data ~ colour + Error(subj/colour),data=table))
```

```
Error: subj
```

	Df	Sum Sq	Mean Sq	F value	Pr(>F)
Residuals	11	0.0053863	0.0004897		

```
Error: subj:colour
```

	Df	Sum Sq	Mean Sq	F value	Pr(>F)
colour	2	0.0037722	0.0018861	95.365	1.447e-11 ***
Residuals	22	0.0004351	0.0000198		

```
---
```

```
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

The statement

`Error(subj/colour)`

is a shorthand for

`Error(subj + subj:colour)`

It tells R to partition the residuals into two error terms, one represents the effects of the differences between subjects (`subj`) and the other is the subjects by colour interaction (`subj:colour`) which is the appropriate error term for testing the effects of `colour` on the RTs.

To understand this procedure we have to remember that while in a fully randomised design we use a common error term to test the effects of all the different factors and their interactions, in a repeated measures design we have to split up this error term into different partitions, some of which we will use to test the effects of our factors and their interactions.

In the case at hand the subjects' error term will not be used to test any effects, it will just be subtracted from the common residuals entry, so that the variability due to differences between subjects doesn't inflate the error term we will use to test for the effects of colours. The subjects' error terms is in this case just the Sum of Squares Between.

The second error term, the subject by colour interaction is the one we want to use to test for the effects of `colour`. This term is what's left from the Sum of Squares Total once you have subtracted the effects due to the subjects and the effects due to the colour, so in the specific design of this experiment it represents just random variability, errors,¹

All this is quite tricky at first. Don't worry, remember as a rule of thumb that in a repeated measures design with R you have to add the `Error()` term to the formula, and this error term is defined as `subjects/your within subj factors`. More examples with Between and Within Subjects factors will be presented in the next sections.

As for the results of this analysis, the F statistics for the colour effect is significant. The RTs for detecting stimuli of these three different colours are different.

12.2.2 Two Within Subjects Factors

The basic principles for running a repeated measures ANOVA with more than one within subject factors are the same as in the case of a within subjects factor only, with just some further complications due to the fact that now we also want to test for interactions between our factors. We'll start straight with an example. Let's say your data look something like the ones in table 12.4, imagining that you have 15 rats exposed to two factors (Alcohol

¹A subject is faster with blue stimuli and gets depressed with black ones, while another subject does just the opposite?! We're not interested in these differences in this experiment, though we might want to test similar effects in other case. So for this time this interaction is just errors.

Table 12.4: Data for the repeated measures ANOVA example

Drug		No-Drug	
Alcohol	No-Alcohol	Alcohol	No-Alcohol
7	6	6	4
5	4	5	2
8	7	7	4
8	8	6	5
6	5	5	3
8	7	7	6
5	5	5	4
7	6	6	5
8	7	6	5
7	6	5	4
9	8	5	4
4	4	3	2
7	7	5	3
7	5	5	0
8	7	6	3

and Drug), with two levels each (administered and not-administered), and the dependent variable is the number of social interaction in a cage with other rats. Here each row represents a subject, you need to reorganise the data so that each row contains a single observation, and the different columns represent:

1. an identifier for the subject
2. the values of the dependent variable (the data that you see in the table)
3. the level of the first factor for each observation
4. the level of the second factor for each observation

The measures you have collected are stored in a text file "rats.txt", in the format of the table above, but with no header whatsoever and no number indicating each subject. The first thing we can do, is to read in these data in R as a number vector:

```
> socialint <- scan ("rats.txt")
```

then we need a column specifying to which subject, each observation of the `socialint` vector belongs to. In addition, we want this new vector to be considered as a factor vector rather than as numerical vector. We can use the `rep()` function to do this:


```
> subj <- as.factor(rep(1:15, each = 4))
```

next we need to specify in a new vector the levels of the first factor for each observation. If we use a character vector to store the levels of the factor, it is not necessary to use the `as.factor()` command, as later, when we will put the all the vectors in a data frame, R will automatically interpret character vectors as factors.

```
> alcohol <- rep( c("A1","No-A1"), 30)
```

we do basically the same for the second factor:

```
> drug <- rep( c("Drug","No-Drug"), 15, each = 2)
```

It's almost done, we need now to put all this vectors in a data frame:

```
> rats <- data.frame (subj, socialint, alcohol, drug)
```

Done! Now the analysis:

```
> summary(aov(socialint ~ alcohol * drug + Error(subj/(alcohol * drug)),  
+           data = rats)
```

the above formula specifies our model for the analysis, we are telling R we want to explain the variable `socialint` by the effects of the factors `alcohol`, `drug` and their interaction, in the case of a repeated measures ANOVA however, we also have to specify the error terms that R will use to calculate the F statistics. The statement

```
Error(subj/(alcohol * drug))
```

is a shorthand for

```
Error(subj + subj:alcohol + subj:drug + subj:alcohol:drug)
```

Again, as in the example with only one within subjects factor, we are telling R to partition the residuals into different error terms. The first is the effect due to differences between subjects. It will not be used to test any effects, it will just be subtracted from the residuals to compute the other error terms. The `subj:alcohol` and the `subj:drug` interactions are the error terms to be used to test the effects of `alcohol` and `drug` respectively, while the `subj:alcohol:drug` interaction will be used to test the interaction between between `alcohol` and `drug`. Below there is the output from the analysis:

```
Error: subj  
          Df Sum Sq Mean Sq F value Pr(>F)  
Residuals 14 69.433   4.960
```

```

Error: subj:alcohol
      Df Sum Sq Mean Sq F value    Pr(>F)
alcohol  1 26.6667 26.6667  42.264 1.399e-05 ***
Residuals 14  8.8333  0.6310
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Error: subj:drug
      Df Sum Sq Mean Sq F value    Pr(>F)
drug    1 60.000  60.000  57.931 2.434e-06 ***
Residuals 14 14.500   1.036
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Error: subj:alcohol:drug
      Df Sum Sq Mean Sq F value    Pr(>F)
alcohol:drug  1 4.2667  4.2667  18.474 0.0007362 ***
Residuals    14 3.2333  0.2310
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

```

As you can see R splits the summary into different sections, based on the partition of the error terms that we have specified. Each effect is then tested against its appropriate error term.

The results tell us that there is a significant effect of both the `alcohol` and the `drug` factors, as well as their interaction.

12.2.3 Two Within Subjects Factors and One Between

Now let's see the case in which we also have a between subjects factor. Suppose we want to run again the experiment on the effects of alcohol and drug on the social interactions in rats, but this time we want to use two different species of rats, the yuppy rats and the kilamany rats, as we have reasons to believe that the kilamany will have different reactions to alcohol and drugs from the yuppy, that is the species that we had tested before. So we manage to gather 8 rats from each species and run our experiment. The results are shown in table 12.5.

the data are in the file "`two_within_one_between.txt`", in each row of this file are recorded the number of social interactions for a rat under the 4 experimental conditions it participated in. We need to get the "one row per observation format":

```

> subj <- rep(1:16,each=4) ##read in the data
> subj <- as.factor(subj)
> alcohol <- rep(c("A1","No-A1"),32)

```

Table 12.5: Data for the repeated measures ANOVA example

Drug		No-Drug	
Alcohol	No-Alcohol	Alcohol	No-Alcohol
Yuppy			
7	6	6	4
5	4	5	2
8	7	7	4
8	8	6	5
6	5	5	3
8	7	7	6
5	5	5	4
7	6	6	5
Kalamani			
8	7	6	5
7	6	5	4
9	8	5	4
4	4	3	2
7	7	5	3
7	5	5	0
8	7	6	3
5	6	4	5

```

> alcohol <- as.factor(alcohol)
> drug <- rep(c("Drug", "No-Drug"), 16, each=2)
> drug <- as.factor(drug)
> group <- rep(c("Yuppy", "Kalamani"), each=32)
> group <- as.factor(group)
> dats <- data.frame(subj, socialint, alcohol, drug, group)

```

now the ANOVA

```

> summary(aov(socialint~alcohol*drug*group + Error(subj/(alcohol*drug)),
+           data=dats))

```

Error: subj

	Df	Sum Sq	Mean Sq	F value	Pr(>F)
group	1	2.250	2.250	0.4615	0.508
Residuals	14	68.250	4.875		

Error: subj:alcohol

	Df	Sum Sq	Mean Sq	F value	Pr(>F)
--	----	--------	---------	---------	--------

```

alcohol      1 22.5625 22.5625 22.7658 0.0002982 ***
alcohol:group 1  0.0625  0.0625  0.0631 0.8053669
Residuals    14 13.8750  0.9911
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

```

```

Error: subj:drug
      Df Sum Sq Mean Sq F value    Pr(>F)
drug      1 60.062  60.062 81.0482 3.376e-07 ***
drug:group 1  5.063   5.063  6.8313 0.02043 *
Residuals 14 10.375   0.741
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

```

```

Error: subj:alcohol:drug
      Df    Sum Sq   Mean Sq   F value   Pr(>F)
alcohol:drug      1      4.00      4.00      16 0.001316 **
alcohol:drug:group 1 5.610e-32 5.610e-32 2.244e-31 1.000000
Residuals         14      3.50      0.25
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

```

Chapter 13

How to Adjust the p-values for Multiple Comparisons

The base R program comes with a number of functions to perform multiple comparisons. Here by “multiple comparisons” I mean both post-hoc, and planned comparisons procedures. Not all the possible procedures are available, and some of them might not be applicable to the object resulting from the specific analysis you’re doing. Additional packages might cover your specific needs.

The `p.adjust` function

The function `p.adjust` comes with the base R program, in the package `stats`, so you don’t need to install anything else on your machine apart from R to get it. This function takes a vector of p-values as an argument, and returns a vector of adjusted p-values according to one of the following methods:

- "holm"
- "hochberg"
- "hommel"
- "bonferroni"
- "BH"
- "BY"
- "fdr"
- "none"

So, if for example you’ve run 3 t-test after a one-way ANOVA with 3 groups, to compare each mean group’s mean with the others, and you want to correct the resulting p-values with the Bonferroni procedure, you can use `p.adjust` as follows:

```
> ps <- c(0.001,0.0092,0.037) #p values from the t-tests
> p.adjust(ps, method="bonferroni")
[1] 0.0030 0.0276 0.1110
```

the last line of the example gives the p-values corrected for the number of comparisons made (3 in our case), using the Bonferroni method. As you can see from the output, the first 2 values would still be significant after the correction, while the last one would not.

You can use any of the procedures listed above, instead of the Bonferroni procedure, by setting it in the `method` option. If you don't set this option at all you get the Holm procedure by default. Look up the Reference Manual for further information on these methods.

The method "none" returns the p-values without any adjustment.

The `mt.rawp2adjp` function

The package `multtest` contains a function very similar to `p.adjust`, and offers some other correction methods.

```
> library(multtest)
> ps <- c(0.001,0.0092,0.037) #p values from the t-tests
> mt.rawp2adjp(ps, proc="Bonferroni")
$adjp
      rawp Bonferroni
[1,] 0.0010      0.0030
[2,] 0.0092      0.0276
[3,] 0.0370      0.1110

$index
[1] 1 2 3
```

or, if you want to see the adjusted p-values with more than one method at once:

```
> library(multtest)
> ps <- c(0.001,0.0092,0.037) #p values from the t-tests
> mt.rawp2adjp(ps, proc="Bonferroni")
$adjp
      rawp Bonferroni      SidakSS      Holm      SidakSD
[1,] 0.0010      0.0030 0.002997001 0.0030 0.002997001
[2,] 0.0092      0.0276 0.027346859 0.0184 0.018315360
[3,] 0.0370      0.1110 0.106943653 0.0370 0.037000000

$index
[1] 1 2 3
```

Chapter 14

R Programming

14.1 Control Structures

If..Else Conditional Execution

It is possible to insert and execute control structures directly from the R interpreter, but for the following examples, I'll assume you're writing the commands to a batch file, and then executing them through the `source()` command.

The general form of conditional execution in R is:

```
if (cond){
  do_this
} else {
  do_this_other
}
```

here's a silly example

```
money <- 1300
if (money > 1200){
  print("good!")
} else {
  print("troubles...")
}
```

it's important that the `else` statement is on the same line where the previous command ends (in the above example that's the closing brace on the fourth line), otherwise the interpreter sees it as unrelated to the previous `if` and will give an error (the `if` statement could also be used by itself, so it would be seen as a complete statement if `else` does not appear on the same line).

It is also possible to execute more than one command upon the fulfilment of a given condition:

```

money <- 900
expenses<-1200
if (money > expenses){
  print("good!")
  shopping<- money-expenses
} else {
  print("troubles...")
  shopping<-NA
}

print("Money available for shopping:")
print(shopping)

```

Finally it is possible to add branches to your control structure with the else if statement:

```

expenses<- 1000
laptop <- 1000
if ((money-expenses) > 1000){
  print("great!! buy new laptop")
  shopping<-(money-expenses)-laptop
} else if ((money-expenses) > 0 && (money-expenses) <= 1000){
  print("no laptop, just shopping and save some")
  shopping<- (money-expenses)/2
} else {
  print("troubles...")
  shopping<-NA
}

print("Money available for shopping:")
print(shopping)

```

ifelse

The `ifelse` function is handy for testing all the elements of a vector on a given condition, the general form is:

```
ifelse(condition, value_if_cond_true, value_if_cond_false)
```

for example, let's say we want to categorise the results of a classroom test, scored from 1 to 10 as "pass" if the score was equal to, or greater than 6 and "fail" if the score was less than 6

```

> score <- c(4,7,6,5,8,6,7,)
> admission <- ifelse(score >= 6, "pass", "fail")
> admission
[1] "fail" "pass" "pass" "fail" "pass" "pass" "pass"

```


so, the first argument of the `ifelse` function, is the condition that we want to test, the second argument is the value that should be returned if the condition is met, and the third argument is the value that should be returned if it is not.

xor

The function `xor` implements the exclusive logical “or” operator, that is, it evaluates to true if exclusively one of two alternative conditions is met, otherwise, it evaluates to false. The latter occurs both, when none of the conditions is met and when both are met simultaneously.

```
> a<-6
> xor(a>5,a>7)
[1] TRUE
> xor(a>5,a>3)
[1] FALSE
```

in the first example, only the first condition ($a > 5$) is met, so the function evaluates to true. In the second example, both conditions are satisfied, but since we’re using `xor` and you can have one thing or the other, but not both together, the function evaluates to false.

```
> a
[1] 4 5 7 3 6 7 5 9 3 6 4 8 3 5 1 9
> a<-a[-which(a>5)]
```

14.2 String Processing

One of the strengths of R, in my opinion, lies in the way it deals with character strings. Certain objects, for example dataframes, allow to mix strings with other data types, subsets of certain objects (again dataframes are an example, but also lists), can be easily given meaningful names and retrieved. This adds much flexibility and ease of use of R compared to other languages (e.g. MATLAB). One aspect that is perhaps less known however, are the powerful string processing functions that R gives you. Once you get to know them you’ll realise you can do all your data analysis in R, without the need to use other languages, like python or perl for pre-processing.

The simplest thing you can do with a string, is counting its characters, which you can do with the `nchar` function:

```
> mystr<-’I love R’
> nchar(my_string)
[1] 8
```

The second thing you can do with strings is extracting parts of them. There are various way to achieve this. Two of the most useful functions are `substr` and `strsplit`.

`substr`, as the name suggests, returns part of a string

```
> substr(my_string, start=1, stop=4)
[1] "I lo"
```

if you want to get a portion of a string from some point in the middle, to the end

```
> substr(my_string, start=5, stop=nchar(my_string))
[1] "ve R"
```

`substr` can be also used to replace parts of a string

```
> substr(my_string,start=1,stop=3)<-'qqq'
> my_string
[1] "qqqove R"
```

14.2.1 Using Regular Expressions

```
> b<-c('the','atheist','theme','therion','thin','jjthe')
> ## match only when pattern appears at the beginning
> grep('^the',b,value=TRUE)
[1] "the"      "theme"    "therion"
> ## match only when pattern appears at the end
> grep('the$',b,value=TRUE)
[1] "the"      "jjthe"
> ## match exactly 'the' not followed or preceded by anything
> grep('^the$',b,value=T)
[1] "the"
> ## match 'the' followed by 'i' or 'm'
> grep('the[i,m]',b,value=TRUE)
[1] "atheist" "theme"
> ## match 'the' followed by anything except 'i'
> grep('the[^i]',b,value=TRUE)
[1] "theme"    "therion"
> ## match 'the' followed by anything
> grep('the.',b,value=TRUE)
[1] "atheist" "theme"    "therion"
> ## match 'the' preceded by anything
> grep('.the',b,value=TRUE)
[1] "atheist" "jjthe"
```

`glob2rx` translates a wildcard pattern, as used in most shells (for example for listing files with the Unix `ls`), in a regular expression, so if you're used to wildcards this comes in handy

```
> glob2rx('the*')
[1] "^the"
> glob2rx('the')
[1] "^the$"
```

14.3 Tips and Tricks

14.3.1 Convert a String into a Command

```
> cmd <- "vec <- c(1,2,3)"
> eval(parse(text=cmd))
```

14.4 Creating Simple R Packages

If you start writing your own functions and you use them often, probably you will soon get tired of sourcing the files containing each function to make them available at each session. There are at least two ways around this problem:

- put all your function files in a directory and write a function that systematically sources them all.
- build a R package

The first solution is rather simple, give the `.R` extension to your R function files and put them in a directory. Although there is not a built-in function to source all the R files present in a directory, the documentation for the `source` function gives an example on how to do it (see `?source`):

```
## If you want to source() a bunch of files, something like
## the following may be useful:
sourceDir <- function(path, trace = TRUE, ...) {
  for (nm in list.files(path, pattern = "\\.[RrSsQq]$")) {
    if(trace) cat(nm,":")
    source(file.path(path, nm), ...)
    if(trace) cat("\n")
  }
}
```

the function sources all the files with the `.R` extension found in the directory indicated by the `path` argument. You can copy this function to a file, let's say

`sourceDir.R`, put it in your HOME directory and source it in your `.First` function in `.Rprofile` (see Chapter 15 for details the `.First` function and the `.Rprofile` file)

```
## This goes in .Rprofile in ~/
.First <- function(){
  source("~/sourceDir.R")
}
```

now each time you call `sourceDir` with a directory as an argument, you will have all the functions defined there available. If you want them available at the beginning of each session, just add a call to `sourceDir` for the directories you want to add in your `.First` function as well. So for example, if your R function files are in the directory `myRfunctions`, add the following to your `.First` function:

```
## This goes in .Rprofile in ~/
.First <- function(){
  source("~/sourceDir.R")
  sourceDir("~/myRfunctions")
}
```

Building a R package requires a bit more work. The detailed documentation for doing this is provided in the “**Writing R Extensions**” manual available at the CRAN website <http://cran.r-project.org/>. That documentation looks at best daunting for a beginner, indeed writing a R package is not trivial, however if all you have is pure R code, and you just want to build a simple package for your own use, the task should not be too difficult to achieve. A very useful document is “**An introduction to the R package mechanism**”, it can be found at the following URL <http://biosun1.harvard.edu/courses/individual/bio271/lectures/L6/Rpkg.pdf>. In the following sections I’ll try to explain how to build a simple R package, much of what I say is drawn from the above cited documents.

14.4.1 The Bare Minimum to Create a Package

The quickest way to get started is to use the function `package.skeleton` to create the first “draft” of your package. Start a R session, make sure that there are not R objects in your session, otherwise they will be bundled in your package

```
> rm(list=ls(all=TRUE))
```

now source all the function files you want to include in your package, for example

```
> source("/home/sam/myFunctions.R")
> source("/home/sam/soundFunctions.R")
```

and then call the `package.skeleton` function with the name you want to give to your package as the argument, for example “mypkg”

```
> package.skeleton("mypkg")
```

this will create a directory called `mypkg` with two sub-directories, “R” containing your code, and “man” containing the documentation files. Furthermore a file called `DESCRIPTION` will be created. If the objects you are packaging include datasets, a `data` directory will also be created. These are the essential elements needed to build a package. The exact content of these files and directories, and how to edit them will be explained later, for the moment I’ll give an overview of the steps required to start using your package. The next step consists of building the package. Start a shell (not a R session), move one directory above the `mypkg` directory we’ve just created and give the command

```
$ R CMD build mypkg
```

to build the package, this will create a tar gzipped file with everything necessary to install the package, the next step to do is indeed the installation. I would recommend installing your own packages in a separate directory from the default package installation directory, let’s say `~/personalRLibrary`, to install the package in this directory, still from the shell call

```
$ R CMD INSTALL -l ~/personalRLibrary nameOfTarFile.tar.gz
```

now from a R session you can call your package with

```
> library(mypkg, lib="~/personalRLibrary")
```

if you want to add permanently your personal R library to the library search path, you can add the following line to the `.First` function in your `.Rprofile`

```
## This goes in .Rprofile in ~/
.First <- function(){
  .libPaths(c(.libPaths()), "~/personalRLibrary/")
}
```

in this way, after starting a new session you’ll be able to load your package without having to specify in which library it is located

```
> library(mypkg)
```

The one described above is a very quick but rough way of creating a package, in order to properly create a R package a number of additional steps, like writing the documentation, and adding examples, need to be followed. Some of these steps will be described in the following sections. Always remember that a very useful thing to do when learning how to build a package is to download some source packages and explore their contents.

Editing the DESCRIPTION file

The DESCRIPTION file follows the Debian control format, and has a key-value pair syntax. The default fields created by `package.skeleton` are pretty much self-explanatory. Other fields that can be added are

- **Depends** If your package depends on a particular version of R, or on other packages, these should be listed here. For example:

```
Depends: R (>= 1.9.0), gtools, gdata, stats
```

- **URL** The URL of a website where you can find out more about the package. For example:

```
URL: http://www.example.com
```

Editing the Documentation

The documentation files reside in the `man` directory of your package. There is one documentation file for each function or data set present in the package. The documentation files are written in a L^AT_EX like format called Rd. `package.skeleton` creates a skeleton of the documentation file, which just needs to be edited, the default fields are pretty much self-explanatory. For a more detailed explanation you can read the “**Writing R Extensions**” manual <http://cran.r-project.org/doc/manuals/R-exts.html>. I’ll give you just a few tips:

It is possible to add additional sections beside the default ones, for example it may be useful to add a “Warnings” section if you have any warnings to give on the use of the function

```
\section{Warning}{Calling this function with arguments foo foo2 can cause ...}
```

In the `seealso` section, you can refer to other functions contained in your package, for example

```
\code{\link{functionFoo}}
```

will automatically add a hyperlink to the documentation for the function `functionFoo`

In the “Examples” section, you write code as if you were writing it in a R script. You can use datasets from your own package, or from the standard R dataset. Keep in mind that the examples should be directly executable by the user, either through copy and paste, or through the `example()` function. When the package is installed, the examples will appear in a directory called `R-ex`, however you do not need to bother about this, the R code for your examples needs to be written within the documentation Rd files.

The documentation requires the presence of one or more standard keywords. One way to get a list of these keywords is to download the tarball with the R sources, after unpacking it, you can find the keywords in a file within the `doc` directory called `KEYWORDS.db`.

Converting Rd files to Other Formats

HTML and L^AT_EX versions of the documentation files are automatically produced in the package installation process, you can find them in the `html` and `latex` directories of your package installation directory, respectively. You can also produce a single pdf or dvi file containing all the documentation using the following command from a shell

```
$ ## produce dvi
$ R CMD Rd2dvi /path/to/your/package/sources/
$ ## produce pdf
$ R CMD Rd2dvi --pdf /path/to/your/package/sources/
```

Adding additional Function or Data Files to the Package

Adding additional function files is quite straightforward, the files contained in the R sub-directory of your package directory are plain R files, so you can just write your functions, drop the files with your functions there, and next time you build the package the new functions will be included. The function `prompt` can be used to build the documentation templates for new functions:

```
> myfun <- function(arg){val <- arg + 3}
> prompt(myfun)
```

this will create a `.Rd` file which you can edit, and drop in the `man` directory.

Datasets can be saved using the `save` function:

```
> mydata <- seq(1,10, .1)
> save(mydata, file='mydata.rda')
```

documentation again can be produced using the `prompt` function

```
> prompt(mydata)
```

Checking the Package

The sanity check for the package can be done by issuing the following command from a shell

```
$ R CMD check /path/to/your/package/sources/
```

Chapter 15

Environment Customisation

Perhaps the best way to customise your R environment is through the use of a `.Rprofile` file, that you put in your `HOME` directory. This is a simple text file that is sourced every time R is started, so you can put in it your own functions, and any operations that you would like R to perform at start-up. Also in this file, you can write two special functions, the `.First` is executed first at the beginning of a session, and the `.Last` is executed at the end of a session. The `.First` function is normally used to initialise the environment setting the desired options. Here's an example of a `.Rprofile` file

```
##This is my .Rprofile in ~/
.First <- function(){
  options(prompt=">>> ", continue="+\t") ##change the prompt
  options(digits=5, length=999)         ##display max 5 digits
}

setwd("~/rwork") ##Start R in this directory
```

for a full listing of the options that can be set see `?options`. You can change temporarily the options, only for the running session, directly from R, for example:

```
> options(digits=9)
```

The `.Rprofile` file in the user's `HOME` is only one of the files that can be used to initialise the environment and set the options. The file that is looked up first by R is the one defined by the `R_PROFILE` environment variable if this variable is set. To verify the value of this variable you can use

```
> Sys.getenv("R_PROFILE")
> system("echo $R_PROFILE")
```

if the variable is unset, R looks for a file called `Rprofile.site` that is in the `etc` sub-directory of your R installation directory. To find out your R installation directory on a Unix-like system you can use


```
> system("echo $R_HOME")
```

The `Rprofile.site` or the file pointed to by the `R_PROFILE` environment variable can be used for system-wide configuration. Note that if the `R_PROFILE` environment variable is set the file pointed to by this variable is used, and if this is not `Rprofile.site`, the latter is ignored.

The `.Rprofile` file in the user's `HOME` can be used for user specific initialisation, and the functions written in this file overwrite, or better “mask” functions with the same name defined in either the file pointed to by the `R_PROFILE` variable or in `Rprofile.site`. Moreover a `.Rprofile` file can be put in any directory, then, if R is started from that directory, this file is sourced, and it masks the definitions given in the user's `HOME` `.Rprofile` file, in this way, it is possible to customise the initialisation for a particular data analysis. Finally, a directory specific initialisation can be given in a `.RData` file, the definitions given here mask also the definitions given in any `.Rprofile` files.

Chapter 16

ESS: Using Emacs Speaks Statistics with R

If you have the `ESS` package installed, you can use Emacs for both editing R source files when working in batch mode, and running a R process from within Emacs. ESS provides an extended set of facilities for both these tasks, among these are syntax highlighting, indentation of code and the ability to work with multiple buffers.

To get syntax highlighting, just use the `.R` extension for naming your file.

To start an R session from within Emacs, press `M-x`, type `R` in the minibuffer, and press `Enter`.

Another nice way of running an R session from inside Emacs is to run first a shell in Emacs (press `M-x` and then type `shell` in the minibuffer), and then calling R from that shell. However, this doesn't involve ESS, so you won't have all the features that ESS adds to the Emacs editing facilities.

For sake of clarity, the use of ESS for editing R source files, and for running a R session will be addressed in two separate sections, however this separation is quite artificial, first because the most proficient use of ESS involves editing a `.R` while running a R session, and second because some of the tips given in one section apply also to the usage of ESS illustrated in the other section. Therefore you're invited to at least skim through both sections, even if you're interested on one usage of ESS only.

16.1 Using ESS for Editing and Debugging R Source Files

If you have a basic knowledge of Emacs you will feel at home. A good way to use ESS is to split the Emacs window horizontally (`C-x 2`) and have a R source file in one buffer and a R process running in the other buffer. You

basically write the R code in the first buffer, and then send it to the R process for evaluation. Here are the shortcuts for sending input to the R process:

C-c C-b Evaluate buffer. This means that all the commands present in the source file will be executed

C-c C-j Evaluate only the current line

C-c C-r Evaluate selected region

There is a set of commands that is equivalent to the above, but moves the cursor to the R process window after the evaluation, that is they ‘Evaluate and go’ (to the other window)

C-c M-b Evaluate buffer and go

C-c M-j Evaluate line and go

C-c M-r Evaluate region and go

To comment/uncomment a region you can use

M-x comment-region Comment region

M-x uncomment-region Uncomment region

M-; Comment/Uncomment region

If you want to switch from a buffer to the other one you can use **C-x o**. Moreover, when you are in a buffer you can make the other window scroll without moving the cursor with **C-M-v**.

Italian Keyboard Mapping Issues

Many Italian keyboards don’t have the braces ‘{’ ’}’ and some other symbols like the tilde ~, however, in Linux, if you’ve chosen an Italian keyboard layout they’re mapped somewhere, and you can access them through a combination of keys, likely combinations are:

```
Shift + AltGr + [   for {
Shift + AltGr + ]   for }
AltGr + ^           for ~
```

if you still can’t find them, try pressing **AltGr** with some keys in the upper part of the keyboard, and see if you can spot them.

In Microsoft Windows the braces ‘{’ ’}’ can be accessed as above, while to access the tilde in most applications, including the R console, you

have to write the ASCII code “126” with some modifier function key (in my laptop that’s `Alt+Fn 126`). In Emacs however you can’t write the tilde in this way, a solution is to write the character in octal code by typing:

```
Ctrl-q 176 RETURN
```

16.2 Using ESS to Interact with a R Process

To start a R process type `M-x R`. `C-p` and `C-n`, or `C-↑` and `C-↓` are for scrolling through the command history.

One thing you need to know, is the ESS “smart underscore” behaviour, that is if you press the underscore once, you get the assignment operator ‘<-’, if you press the underscore twice you get a literal underscore ‘_’. This shortcut for the assignment operator can be very annoying if you use the underscore often in variable names, to turn-off this smart behaviour you need to put the following line somewhere in your `.emacs` file, but before ESS is loaded:

```
(ess-toggle-underscore nil)
```

When you are running R inside Emacs, if the cursor is not positioned at the current command line, and you try to retrieve some command from the command history with `C-p` or `C-↑`, Emacs will complain saying that you are “not at command line”. To get at the command line without using the mouse, press `M-Shift-⏪`.

If you set the cursor at a previous command, and then press `Enter`, that command will be evaluated again.

There is a quick way to source a file, `C-c C-l filename` will load and source your file.

`C-c C-q` is another way of quitting R in ESS mode, but I guess `C-d` remains the fastest one.

Chapter 17

Using Sweave to Write Documents with R and \LaTeX

17.1 What's this Sweave?

Sweave provides a great way of writing documents or reports that contain both text and R objects, namely figures, syntax and raw or nicely formatted output produced by R. Sweave is also a way to automate the production of documents. Usually when you write a document with some statistical content, you first write the text, and then add the graphics produced by some statistical software or spreadsheet application. Of course this process is time consuming, not to mention the fact that often you have to manually fill in tables with the statistics produced by the application you're using. Even worse, if you've already prepared your report, but then you have to change something in the statistical analyses, for example you have to add or drop a subject, the entire process must be repeated again from scratch. Sweave is aimed at easing this process, it works like this: you write a single source file which contains your text, written with the \LaTeX markup language, and embed in it the R syntax needed for producing the figures, and tables to appear in your document. You process this file with Sweave through R, and you get a plain \LaTeX source file which you can run with, well \LaTeX of course, to get your nicely formatted output. If your data happen to change for any reasons, you don't have to start again from scratch, you can just run your old Sweave source file on your new dataset, and everything will be updated automatically.

If you already know \LaTeX and R, learning to use Sweave will be easy, the additional syntax required by Sweave to integrate R and \LaTeX source code is minimal. As a caveat, I have to say that Sweave is still in its infancy, so, for the moment probably you won't automatically get all the tables, with

complicated layouts that you'd like to add to your document. In my opinion, however, Sweave already does a great job, and it's well worth using.

17.2 Usage

Recent versions of R come with Sweave already in the base system (in the package `utils`), so you don't need to install it separately. Of course you need \LaTeX installed to produce the document later, Sweave only outputs a \LaTeX source file and all the graphics needed for your document.

You start the Sweave source file as a normal \LaTeX document with the usual preamble, if you name the file with the `.Rnw` extension, Emacs should recognise it and highlight the syntax for you. Then, at the point where you want to embed a chunk of R syntax, you add the following tag:

```
<<>>=
```

after you've finished with the R syntax chunk, you need to add the following tag to start another piece of text written with \LaTeX :

```
@
```

in this way you can alternate chunks of R code with pieces of \LaTeX syntax. If you forget to add the `@` tag before a piece of \LaTeX syntax Sweave will complain and abort the process, if you instead make a mistake with \LaTeX syntax, Sweave won't complain and will normally produce the `.tex` file, however the compilation with \LaTeX won't work.

There are different options that determine which R objects will appear in the final document. If you want both the R code, and its output to appear in the document, just use the `<<>>=` empty tag. They will both be inserted in the \LaTeX file in a redefined `verbatim` environment. Setting the option `echo=FALSE` lines of R code are not included in the document, while with the option `results=hide` the output of the R code won't appear in the document. Therefore if you want to run some R code, but neither the code nor its output should appear in the document, use:

```
<<echo=FALSE, results=hide>>=
```

Another option will suppress all the output, except the figures:

```
<<echo=FALSE, fig=TRUE>>=
```

Sweave will automatically put a `\includegraphics{}` command for a figure.

Finally, if you want to use some utilities, like `xtable`, that automatically produce \LaTeX objects from R objects, you will want to use the following options to tell Sweave not to put the R output in a `verbatim` environment:

```
<<echo=FALSE, results=tex>>=
```

Chapter 18

Sound Processing

18.1 Libraries for Sound Analysis and Signal Processing

seewave

seewave provides functions for analysing, manipulating, displaying, editing and synthesising time waves (particularly sound). This package processes time analysis (oscillograms and envelopes), spectral content, resonance quality factor, cross correlation and autocorrelation, zero-crossing, dominant frequency, 2D and 3D spectrograms.

<http://cran.r-project.org/src/contrib/Descriptions/seewave.html>

sound

Basic functions for dealing with wav files and sound samples.

<http://cran.r-project.org/src/contrib/Descriptions/sound.html>

tuneR

Collection of tools to analyse music, handling wave files, transcription.

<http://cran.r-project.org/src/contrib/Descriptions/tuneR.html>

signal

A set of generally Matlab/Octave-compatible signal processing functions. Includes filter generation utilities, filtering functions, re-sampling routines, and visualisation of filter models. It also includes interpolation functions and some Matlab compatibility functions.

<http://cran.r-project.org/src/contrib/Descriptions/signal.html>

Appendix A

Partial List of Packages by Category

A.1 Graphics Packages

lattice

<http://cran.r-project.org/src/contrib/Descriptions/lattice.html>

rgl

<http://cran.r-project.org/src/contrib/Descriptions/rgl.html>

tkrplot

<http://cran.r-project.org/src/contrib/Descriptions/tkrplot.html>

iplots

<http://cran.r-project.org/src/contrib/Descriptions/iplots.html>

rpanel

rpanel provides a set of functions to build simple GUI controls for R functions. These are built on the tcltk package. Uses could include changing a parameter on a graph by animating it with a slider or a "doublebutton", up to more sophisticated control panels.

<http://cran.r-project.org/src/contrib/Descriptions/rpanel.html>

A.2 GUI Packages

Rcmdr

<http://cran.r-project.org/src/contrib/Descriptions/Rcmdr.html>

JGR

<http://cran.r-project.org/src/contrib/Descriptions/JGR.html>

pmg

Simple GUI for R using gWidgets <http://cran.r-project.org/src/contrib/Descriptions/pmg.html>
<http://www.math.csi.cuny.edu/pmg>

Appendix B

Miscellaneous commands

Data

`head(dats)` print the first part of the object `dats`

`tail(dats)` print the last part of the object `dats`

Help

`?foo` find help on command `foo`

Objects

`class(foo)` get the class of object `"foo"`

`ls()` `objects()` list objects present in the workspace

Organise a session

`dir()` or `list.files()` list the files in the current directory

`getwd()` get the current directory

`library(foo)` load library `foo`

`library()` list all available packages

`q()` or `quit()` quit from current session

`require(foo)` require the library `foo`, use in scripts

`setwd("home/foo")` set the working directory

`system("foo")` execute the system command `foo` as if it were from the shell

R administration

`library()` list all installed packages

`R.version` info on R version and the platform it is running on

Syntax

`#` starts a comment

`mydata$foo` refer to the variable `foo` in the dataframe `mydata`

`:` interaction operator for model formulae, `a:b` is the interaction between `a` and `b`

`*` crossing operator for model formulae, `a*b = a+b+a:b`

Appendix C

Other Manuals and Sources of Information on R

- **R Project Homepage**
<http://www.r-project.org/>
- **CRAN**
<http://cran.r-project.org/>
- **R mailing lists**
<http://www.r-project.org/mail.html>
- **Notes on the use of R for psychology experiments and questionnaires** by J. Baron and Y. Li
<http://www.psych.upenn.edu/~baron/rpsych.pdf>
- **R manual for biometry** by K.J. Hoff
http://www.bioinf.uni-hannover.de/teaching/Hoff/RMANUAL_ENGLISH.PDF
- **Introductory R course at the University of Lancaster**
http://www.cas.lancs.ac.uk/short_courses/intro_r.html
- **Generalised linear models course with R at the University of Lancaster**
http://www.cas.lancs.ac.uk/short_courses/gen_models.html
- **Simple R** by J. Verzani
<http://www.math.csi.cuny.edu/Statistics/R/simpleR>
- **Simplified manual of some R functions for the Stem and Tendril project**
<http://wiener.math.csi.cuny.edu/st/stRmanual>
- **A R wiki**
http://www.sciviews.org/_rgui/wiki

- **Using R for psychological research: A very simple guide to a very elegant package** by William Revelle
<http://personality-project.org/r/>
- **Practical Regression and Anova Using R** by Julian J. Faraway
<http://www.stat.lsa.umich.edu/~faraway/book/practical.pdf>
- **Jonathan Baron's R page**
<http://finzi.psych.upenn.edu/>
- **Yuelin Li's R page**
<http://idecide.mskcc.org/stats/>
- **Deepayan Sarkar Page**
<http://www.stat.wisc.edu/~deepayan/>
- **Vincent Zoonekynd's R page**
http://zoonek2.free.fr/UNIX/48_R/all.html
- **Christian Hoffman's R-Howto**
<http://www.wsl.ch/staff/christian.hoffmann/R-Howto.ehtml>
- **P.M.E. Altham's page on multivariate analysis with R notes**
<http://www.statslab.cam.ac.uk/~pat/>
- **Ecological models and data in R** by Ben Bolker
<http://www.zoo.ufl.edu/bolker/emd/book/>
- **Patrick Burns' R page**
<http://www.burns-stat.com/>

In Italian

- **Laboratorio di statistica con R** Università Ca' Foscari di Venezia, Dipartimento di Statistica
<http://www.dst.unive.it/~laboratorio/download-materiale.html>
- **Luca Scrucca R's page**
<http://www.stat.unipg.it/~luca/>
- **Roberto Boggiani R's page**
<http://digilander.libero.it/robicox/>
- **Tecniche multivariate** di Corrado Caudek, Università di Trieste, Dipartimento di Psicologia
<http://www.psico.univ.trieste.it/~caudek/multivariata/>

Appendix D

Other Useful Statistics Resources




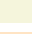
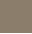
- **Statistics short courses at the University of Lancaster**
http://www.cas.lancs.ac.uk/short_courses/outline_methodology.html#mfs
- **David Howell's statistics page**
<http://www.uvm.edu/~dhowell/StatPages/StatHomePage.html>
- **Weaver's stats page**
<http://www.angelfire.com/wv/bwhomedir/stats.html>

In Italian

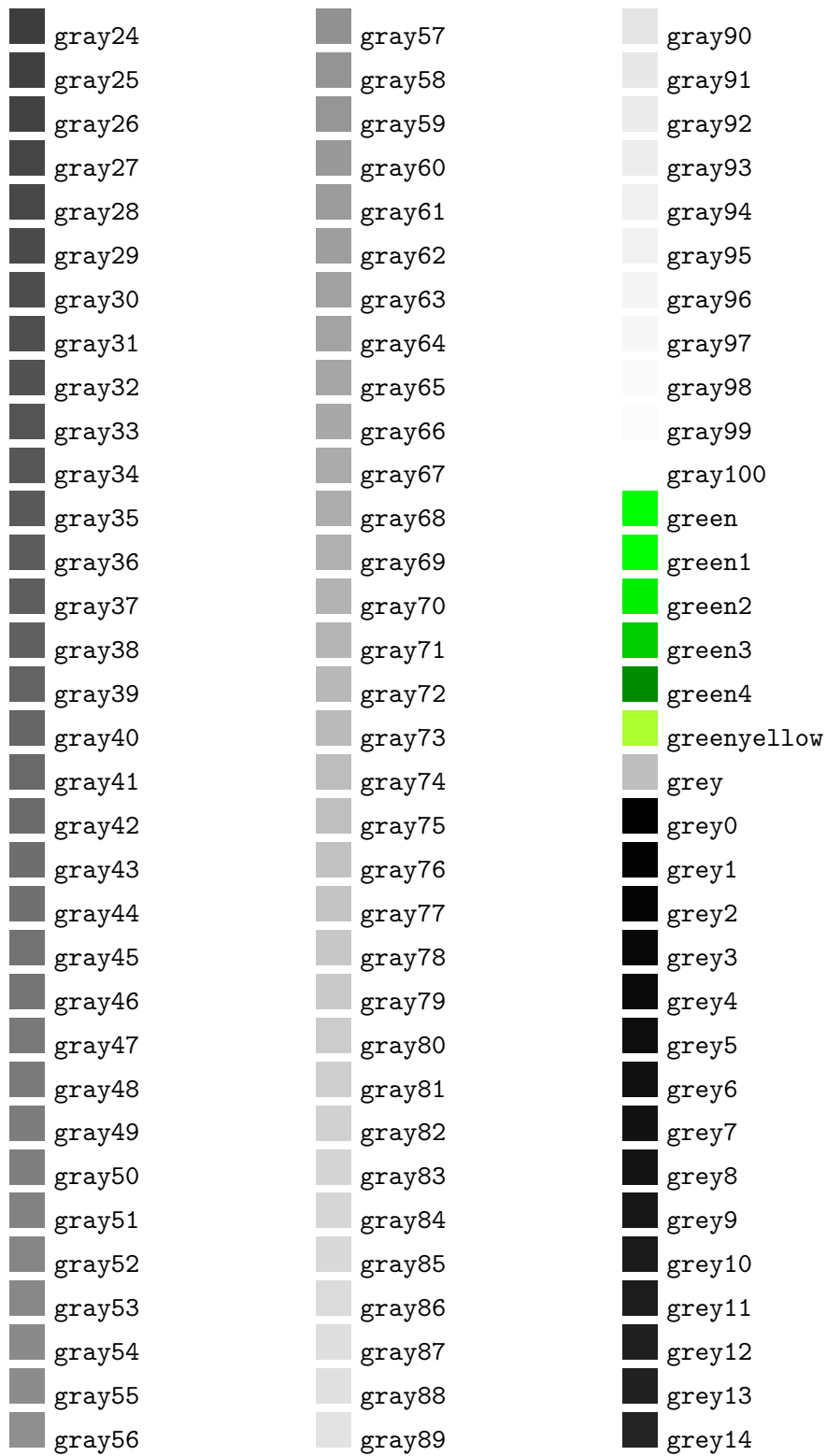
- **Statistica univariata e bivariata parametrica e non-parametrica per le discipline ambientali e biologiche** di Lamberto Soliani
<http://www.dsa.unipr.it/soliani/soliani.html>

Appendix E

Full Colours Table

	white		blanchedalmond		chartreuse3
	aliceblue		blue		chartreuse4
	antiquewhite		blue1		chocolate
	antiquewhite1		blue2		chocolate1
	antiquewhite2		blue3		chocolate2
	antiquewhite3		blue4		chocolate3
	antiquewhite4		blueviolet		chocolate4
	aquamarine		brown		coral
	aquamarine1		brown1		coral1
	aquamarine2		brown2		coral2
	aquamarine3		brown3		coral3
	aquamarine4		brown4		coral4
	azure		burlywood		cornflowerblue
	azure1		burlywood1		cornsilk
	azure2		burlywood2		cornsilk1
	azure3		burlywood3		cornsilk2
	azure4		burlywood4		cornsilk3
	beige		cadetblue		cornsilk4
	bisque		cadetblue1		cyan
	bisque1		cadetblue2		cyan1
	bisque2		cadetblue3		cyan2
	bisque3		cadetblue4		cyan3
	bisque4		chartreuse		cyan4
	black		chartreuse1		darkblue
			chartreuse2		darkcyan

	darkgoldenrod		darkslategray2		gold2
	darkgoldenrod1		darkslategray3		gold3
	darkgoldenrod2		darkslategray4		gold4
	darkgoldenrod3		darkslategrey		goldenrod
	darkgoldenrod4		darkturquoise		goldenrod1
	darkgray		darkviolet		goldenrod2
	darkgreen		deeppink		goldenrod3
	darkgrey		deeppink1		goldenrod4
	darkkhaki		deeppink2		gray
	darkmagenta		deeppink3		gray0
	darkolivegreen		deeppink4		gray1
	darkolivegreen1		deepskyblue		gray2
	darkolivegreen2		deepskyblue1		gray3
	darkolivegreen3		deepskyblue2		gray4
	darkolivegreen4		deepskyblue3		gray5
	darkorange		deepskyblue4		gray6
	darkorange1		dimgray		gray7
	darkorange2		dimgrey		gray8
	darkorange3		dodgerblue		gray9
	darkorange4		dodgerblue1		gray10
	darkorchid		dodgerblue2		gray11
	darkorchid1		dodgerblue3		gray12
	darkorchid2		dodgerblue4		gray13
	darkorchid3		firebrick		gray14
	darkorchid4		firebrick1		gray15
	darkred		firebrick2		gray16
	darksalmon		firebrick3		gray17
	darkseagreen		firebrick4		gray18
	darkseagreen1		floralwhite		gray19
	darkseagreen2		forestgreen		gray20
	darkseagreen3		gainsboro		gray21
	darkseagreen4		ghostwhite		gray22
	darkslateblue		gold		gray23
	darkslategray		gold1		
	darkslategray1				











Appendix F

GNU Free Documentation License

Version 1.2, November 2002

Copyright ©2000,2001,2002 Free Software Foundation, Inc.

51 Franklin St, Fifth Floor, Boston, MA 02110-1301 USA

Everyone is permitted to copy and distribute verbatim copies of this license document, but changing it is not allowed.

Preamble

The purpose of this License is to make a manual, textbook, or other functional and useful document "free" in the sense of freedom: to assure everyone the effective freedom to copy and redistribute it, with or without modifying it, either commercially or noncommercially. Secondly, this License preserves for the author and publisher a way to get credit for their work, while not being considered responsible for modifications made by others.

This License is a kind of "copyleft", which means that derivative works of the document must themselves be free in the same sense. It complements the GNU General Public License, which is a copyleft license designed for free software.

We have designed this License in order to use it for manuals for free software, because free software needs free documentation: a free program should come with manuals providing the same freedoms that the software does. But this License is not limited to software manuals; it can be used for any textual work, regardless of subject matter or whether it is published as a printed book. We recommend this License principally for works whose purpose is instruction or reference.

1. APPLICABILITY AND DEFINITIONS

This License applies to any manual or other work, in any medium, that contains a notice placed by the copyright holder saying it can be distributed under the terms of this License. Such a notice grants a world-wide, royalty-free license, unlimited in duration, to use that work under the conditions stated herein. The "**Document**", below, refers to any such manual or work. Any member of the public is a licensee, and is addressed as "**you**". You accept the license if you copy, modify or distribute the work in a way requiring permission under copyright law.

A "**Modified Version**" of the Document means any work containing the Document or a portion of it, either copied verbatim, or with modifications and/or translated into another language.

A "**Secondary Section**" is a named appendix or a front-matter section of the Document that deals exclusively with the relationship of the publishers or authors of the Document to the Document's overall subject (or to related matters) and contains nothing that could fall directly within that overall subject. (Thus, if the Document is in part a textbook of mathematics, a Secondary Section may not explain any mathematics.) The relationship could be a matter of historical connection with the subject or with related matters, or of legal, commercial, philosophical, ethical or political position regarding them.

The "**Invariant Sections**" are certain Secondary Sections whose titles are designated, as being those of Invariant Sections, in the notice that says that the Document is released under this License. If a section does not fit the above definition of Secondary then it is not allowed to be designated as Invariant. The Document may contain zero Invariant Sections. If the Document does not identify any Invariant Sections then there are none.

The "**Cover Texts**" are certain short passages of text that are listed, as Front-Cover Texts or Back-Cover Texts, in the notice that says that the Document is released under this License. A Front-Cover Text may be at most 5 words, and a Back-Cover Text may be at most 25 words.

A "**Transparent**" copy of the Document means a machine-readable copy, represented in a format whose specification is available to the general public, that is suitable for revising the document straightforwardly with generic text editors or (for images composed of pixels) generic paint programs or (for drawings) some widely available drawing editor, and that is suitable for input to text formatters or for automatic translation to a variety of formats suitable for input to text formatters. A copy made in an otherwise Transparent file format whose markup, or absence of markup, has been arranged to thwart or discourage subsequent modification by readers is not Transparent. An image format is not Transparent if used for any substantial amount of text. A copy that is not "Transparent" is called "**Opaque**".

Examples of suitable formats for Transparent copies include plain ASCII without markup, Texinfo input format, LaTeX input format, SGML or XML using a publicly available DTD, and standard-conforming simple HTML,

PostScript or PDF designed for human modification. Examples of transparent image formats include PNG, XCF and JPG. Opaque formats include proprietary formats that can be read and edited only by proprietary word processors, SGML or XML for which the DTD and/or processing tools are not generally available, and the machine-generated HTML, PostScript or PDF produced by some word processors for output purposes only.

The "**Title Page**" means, for a printed book, the title page itself, plus such following pages as are needed to hold, legibly, the material this License requires to appear in the title page. For works in formats which do not have any title page as such, "Title Page" means the text near the most prominent appearance of the work's title, preceding the beginning of the body of the text.

A section "**Entitled XYZ**" means a named subunit of the Document whose title either is precisely XYZ or contains XYZ in parentheses following text that translates XYZ in another language. (Here XYZ stands for a specific section name mentioned below, such as "**Acknowledgements**", "**Dedications**", "**Endorsements**", or "**History**".) To "**Preserve the Title**" of such a section when you modify the Document means that it remains a section "Entitled XYZ" according to this definition.

The Document may include Warranty Disclaimers next to the notice which states that this License applies to the Document. These Warranty Disclaimers are considered to be included by reference in this License, but only as regards disclaiming warranties: any other implication that these Warranty Disclaimers may have is void and has no effect on the meaning of this License.

2. VERBATIM COPYING

You may copy and distribute the Document in any medium, either commercially or noncommercially, provided that this License, the copyright notices, and the license notice saying this License applies to the Document are reproduced in all copies, and that you add no other conditions whatsoever to those of this License. You may not use technical measures to obstruct or control the reading or further copying of the copies you make or distribute. However, you may accept compensation in exchange for copies. If you distribute a large enough number of copies you must also follow the conditions in section 3.

You may also lend copies, under the same conditions stated above, and you may publicly display copies.

3. COPYING IN QUANTITY

If you publish printed copies (or copies in media that commonly have printed covers) of the Document, numbering more than 100, and the Document's license notice requires Cover Texts, you must enclose the copies in

covers that carry, clearly and legibly, all these Cover Texts: Front-Cover Texts on the front cover, and Back-Cover Texts on the back cover. Both covers must also clearly and legibly identify you as the publisher of these copies. The front cover must present the full title with all words of the title equally prominent and visible. You may add other material on the covers in addition. Copying with changes limited to the covers, as long as they preserve the title of the Document and satisfy these conditions, can be treated as verbatim copying in other respects.

If the required texts for either cover are too voluminous to fit legibly, you should put the first ones listed (as many as fit reasonably) on the actual cover, and continue the rest onto adjacent pages.

If you publish or distribute Opaque copies of the Document numbering more than 100, you must either include a machine-readable Transparent copy along with each Opaque copy, or state in or with each Opaque copy a computer-network location from which the general network-using public has access to download using public-standard network protocols a complete Transparent copy of the Document, free of added material. If you use the latter option, you must take reasonably prudent steps, when you begin distribution of Opaque copies in quantity, to ensure that this Transparent copy will remain thus accessible at the stated location until at least one year after the last time you distribute an Opaque copy (directly or through your agents or retailers) of that edition to the public.

It is requested, but not required, that you contact the authors of the Document well before redistributing any large number of copies, to give them a chance to provide you with an updated version of the Document.

4. MODIFICATIONS

You may copy and distribute a Modified Version of the Document under the conditions of sections 2 and 3 above, provided that you release the Modified Version under precisely this License, with the Modified Version filling the role of the Document, thus licensing distribution and modification of the Modified Version to whoever possesses a copy of it. In addition, you must do these things in the Modified Version:

- A. Use in the Title Page (and on the covers, if any) a title distinct from that of the Document, and from those of previous versions (which should, if there were any, be listed in the History section of the Document). You may use the same title as a previous version if the original publisher of that version gives permission.
- B. List on the Title Page, as authors, one or more persons or entities responsible for authorship of the modifications in the Modified Version, together with at least five of the principal authors of the Document (all

of its principal authors, if it has fewer than five), unless they release you from this requirement.

- C. State on the Title page the name of the publisher of the Modified Version, as the publisher.
- D. Preserve all the copyright notices of the Document.
- E. Add an appropriate copyright notice for your modifications adjacent to the other copyright notices.
- F. Include, immediately after the copyright notices, a license notice giving the public permission to use the Modified Version under the terms of this License, in the form shown in the Addendum below.
- G. Preserve in that license notice the full lists of Invariant Sections and required Cover Texts given in the Document's license notice.
- H. Include an unaltered copy of this License.
- I. Preserve the section Entitled "History", Preserve its Title, and add to it an item stating at least the title, year, new authors, and publisher of the Modified Version as given on the Title Page. If there is no section Entitled "History" in the Document, create one stating the title, year, authors, and publisher of the Document as given on its Title Page, then add an item describing the Modified Version as stated in the previous sentence.
- J. Preserve the network location, if any, given in the Document for public access to a Transparent copy of the Document, and likewise the network locations given in the Document for previous versions it was based on. These may be placed in the "History" section. You may omit a network location for a work that was published at least four years before the Document itself, or if the original publisher of the version it refers to gives permission.
- K. For any section Entitled "Acknowledgements" or "Dedications", Preserve the Title of the section, and preserve in the section all the substance and tone of each of the contributor acknowledgements and/or dedications given therein.
- L. Preserve all the Invariant Sections of the Document, unaltered in their text and in their titles. Section numbers or the equivalent are not considered part of the section titles.
- M. Delete any section Entitled "Endorsements". Such a section may not be included in the Modified Version.

- N. Do not retitle any existing section to be Entitled "Endorsements" or to conflict in title with any Invariant Section.
- O. Preserve any Warranty Disclaimers.

If the Modified Version includes new front-matter sections or appendices that qualify as Secondary Sections and contain no material copied from the Document, you may at your option designate some or all of these sections as invariant. To do this, add their titles to the list of Invariant Sections in the Modified Version's license notice. These titles must be distinct from any other section titles.

You may add a section Entitled "Endorsements", provided it contains nothing but endorsements of your Modified Version by various parties—for example, statements of peer review or that the text has been approved by an organization as the authoritative definition of a standard.

You may add a passage of up to five words as a Front-Cover Text, and a passage of up to 25 words as a Back-Cover Text, to the end of the list of Cover Texts in the Modified Version. Only one passage of Front-Cover Text and one of Back-Cover Text may be added by (or through arrangements made by) any one entity. If the Document already includes a cover text for the same cover, previously added by you or by arrangement made by the same entity you are acting on behalf of, you may not add another; but you may replace the old one, on explicit permission from the previous publisher that added the old one.

The author(s) and publisher(s) of the Document do not by this License give permission to use their names for publicity for or to assert or imply endorsement of any Modified Version.

5. COMBINING DOCUMENTS

You may combine the Document with other documents released under this License, under the terms defined in section 4 above for modified versions, provided that you include in the combination all of the Invariant Sections of all of the original documents, unmodified, and list them all as Invariant Sections of your combined work in its license notice, and that you preserve all their Warranty Disclaimers.

The combined work need only contain one copy of this License, and multiple identical Invariant Sections may be replaced with a single copy. If there are multiple Invariant Sections with the same name but different contents, make the title of each such section unique by adding at the end of it, in parentheses, the name of the original author or publisher of that section if known, or else a unique number. Make the same adjustment to the section titles in the list of Invariant Sections in the license notice of the combined work.

In the combination, you must combine any sections Entitled "History" in the various original documents, forming one section Entitled "History"; likewise combine any sections Entitled "Acknowledgements", and any sections Entitled "Dedications". You must delete all sections Entitled "Endorsements".

6. COLLECTIONS OF DOCUMENTS

You may make a collection consisting of the Document and other documents released under this License, and replace the individual copies of this License in the various documents with a single copy that is included in the collection, provided that you follow the rules of this License for verbatim copying of each of the documents in all other respects.

You may extract a single document from such a collection, and distribute it individually under this License, provided you insert a copy of this License into the extracted document, and follow this License in all other respects regarding verbatim copying of that document.

7. AGGREGATION WITH INDEPENDENT WORKS

A compilation of the Document or its derivatives with other separate and independent documents or works, in or on a volume of a storage or distribution medium, is called an "aggregate" if the copyright resulting from the compilation is not used to limit the legal rights of the compilation's users beyond what the individual works permit. When the Document is included in an aggregate, this License does not apply to the other works in the aggregate which are not themselves derivative works of the Document.

If the Cover Text requirement of section 3 is applicable to these copies of the Document, then if the Document is less than one half of the entire aggregate, the Document's Cover Texts may be placed on covers that bracket the Document within the aggregate, or the electronic equivalent of covers if the Document is in electronic form. Otherwise they must appear on printed covers that bracket the whole aggregate.

8. TRANSLATION

Translation is considered a kind of modification, so you may distribute translations of the Document under the terms of section 4. Replacing Invariant Sections with translations requires special permission from their copyright holders, but you may include translations of some or all Invariant Sections in addition to the original versions of these Invariant Sections. You may include a translation of this License, and all the license notices in the Document, and any Warranty Disclaimers, provided that you also include the original English version of this License and the original versions of those

notices and disclaimers. In case of a disagreement between the translation and the original version of this License or a notice or disclaimer, the original version will prevail.

If a section in the Document is Entitled "Acknowledgements", "Dedications", or "History", the requirement (section 4) to Preserve its Title (section 1) will typically require changing the actual title.

9. TERMINATION

You may not copy, modify, sublicense, or distribute the Document except as expressly provided for under this License. Any other attempt to copy, modify, sublicense or distribute the Document is void, and will automatically terminate your rights under this License. However, parties who have received copies, or rights, from you under this License will not have their licenses terminated so long as such parties remain in full compliance.

10. FUTURE REVISIONS OF THIS LICENSE

The Free Software Foundation may publish new, revised versions of the GNU Free Documentation License from time to time. Such new versions will be similar in spirit to the present version, but may differ in detail to address new problems or concerns. See <http://www.gnu.org/copyleft/>.

Each version of the License is given a distinguishing version number. If the Document specifies that a particular numbered version of this License "or any later version" applies to it, you have the option of following the terms and conditions either of that specified version or of any later version that has been published (not as a draft) by the Free Software Foundation. If the Document does not specify a version number of this License, you may choose any version ever published (not as a draft) by the Free Software Foundation.

ADDENDUM: How to use this License for your documents

To use this License in a document you have written, include a copy of the License in the document and put the following copyright and license notices just after the title page:

Copyright ©YEAR YOUR NAME. Permission is granted to copy, distribute and/or modify this document under the terms of the GNU Free Documentation License, Version 1.2 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".

If you have Invariant Sections, Front-Cover Texts and Back-Cover Texts, replace the "with...Texts." line with this:

with the Invariant Sections being LIST THEIR TITLES, with the Front-Cover Texts being LIST, and with the Back-Cover Texts being LIST.

If you have Invariant Sections without Cover Texts, or some other combination of the three, merge those two alternatives to suit the situation.

If your document contains nontrivial examples of program code, we recommend releasing these examples in parallel under your choice of free software license, such as the GNU General Public License, to permit their use in free software.

References

- [1] J. Baron and Y. Li. *Notes on the use of R for psychology experiments and questionnaires*. 2003.
- [2] R. A. Becker and W. S. Cleveland. *S-PLUS Trellis Graphics User's Manual*. Seattle: MathSoft, Inc., Murray Hill: Bell Labs, 2002.
- [3] R. A. Becker, W. S. Cleveland, and M.-J. Shyu. The visual design and control of Trellis display. *Journal of Computational and Graphical Statistics*, 5:123–155, 1996.
- [4] R. A. Becker, W. S. Cleveland, M.-J. Shyu, and S. P. Kaluzny. *A Tour of Trellis Graphics*. 1996. Available: <http://cm.bell-labs.com/stat/doc/trellis.tour.col.ps>.
- [5] W. S. Cleveland and M. Fuentes. Trellis display: Modeling data from designed experiments. Technical report, Bell Labs, 1997.
- [6] P. Dalgaard. *Introductory Statistics with R*. Springer, 2002.
- [7] J. J. Faraway. *Practical Regression and Anova using R*. 2002.
- [8] R. M. Heiberger and B. Holland. *Statistical Analysis and Data Display. An Intermediate Course with Examples in S-Plus, R, and SAS*. Springer, 2004.
- [9] K. J. Hoff. *R-Manual for Biometry*. 2005. Unpublished Bachelor Thesis at the Teaching Unit Bioinformatics, University of Hannover.
- [10] D. C. Howell. *Statistical Methods for Psychology*. Duxbury, Pacific Grove CA (USA), fifth edition, 2002.
- [11] A. Krause and M. Olson. *The Basics of S-Plus*. Springer, third edition, 2002.
- [12] F. Leisch. *Sweave User Manual*, 2005.
- [13] A. M. Mineo. *Una Guida all'utilizzo dell'Ambiente Statistico R*. 2003.

- [14] P. Murrell. R Lattice graphics. In K. Hornik and F. Leisch, editors, *Proceedings of the 2nd International Workshop on Distributed Statistical Computing, March 15-17, 2001, Technische Universität Wien, Vienna, Austria*, 2001. ISSN 1609-395X.
- [15] A. J. Rossini, R. M. Heiberger, K. Hornik, M. Maechler, R. A. Sparapani, and S. J. Eglen. *ESS—Emacs Speaks Statistics, for version 5.2.11*. The ESS Developers, 2005.
- [16] D. Sarkar. Lattice, an implementation of trellis graphics in r. *R News*, 2(2):19–22, 2002.
- [17] D. Sarkar. Some notes on lattice. In K. Hornik and F. Leisch, editors, *Proceedings of the 3rd International Workshop on Distributed Statistical Computing, March 20-22, 2003, Technische Universität Wien, Vienna, Austria*, 2003. ISSN 1609-395X.
- [18] L. Soliani. *Statistica univariata e bivariata parametrica e non-parametrica per le discipline ambientali e biologiche*. 2005.
- [19] R. Stallman. *GNU Emacs Manual, for version 21.3*. Free Software Foundation, 14th edition, 2002.
- [20] The R Development Core Team. *R Data Import/Export*. Version 2.1.0 (2005-04-18).
- [21] The R Development Core Team. *R Documentation*, 2005.
- [22] W. N. Venables and B. D. Ripley. *S Programming*. Springer, 2000.
- [23] W. N. Venables and B. D. Ripley. *Modern Applied Statistics with S*. Springer, Fourth edition, 2002.
- [24] W. N. Venables, D. M. Smith, and the R Development Core Team. *An Introduction to R*. The R Development Core Team, 2005. Version 2.1.0.
- [25] J. Verzani. *simpleR – Using R for Introductory Statistics*. 2002. Version 0.4 Available: <http://www.math.csi.cuny.edu/Statistics/R/simpleR>.

Index

- ANOVA, **97**
 - one-way, 97
 - one-way unequal sample sizes, 99
 - repeated measures, 100
- assignment, 12
- attach(), 21
- barplot, **42**
 - barplot2, 45
 - with error bars, 45
- bitmap, 72
- boxplot, **47**, 102

- cat, 37
- colour, **69**
 - list, 70
- control structures, **111**

- dataframes, **19**
 - create, 20
 - naming elements of, 21
 - subsetting, 20, 22
- detach(), 21

- Emacs Speaks Statistics, **122**

- graphics
 - devices, 72
 - export, 72
 - with \LaTeX symbols, 59

- head, **32**
- histogram, **49**

- if...else, **111**
- ifelse, **112**
- interaction plot, **51**

- JPEG, 72

- latex
 - in graphics, 59
 - Sweave, 125
- legend, 58
- Levene test, **92**
- lists, **18**
 - naming elements of, 18
 - subsetting, 18
- logical operators, 14
- lty, **55**

- matplot, **42**
- matrix, **15**
 - operations, 17
 - subsetting, 16
 - transpose, 17

- normal distribution, **88**

- oneway.test, **98**

- par, **54**
 - mfcoll, 60
 - mfg, 61
 - mfrow, 60
- pch, **57**
- pdf, 72
- plot
 - add text, 59
- png, 72
- postscript, 72

- read.table, 34
- relational operators, 14
- rep(), **15**, 98

scale(), **38**
scan, 34
scan(), 97
seq(), **14**
sequences, **14**
stack, **29**
stripchart, **49**
Sweave, **125**
symbols, **57**

t-test, **90**
 one-sample, 90
 two samples, 91
tail, **32**

unstack, **29**

vectors, 12
 indexing, 13
 naming elements of, 14
 of characters, 13

write.table, 36

xor, **113**