

An Introduction to R

Notes on R: A Programming Environment for Data Analysis and Graphics
Version 3.1.2 (2014-10-31)

W. N. Venables, D. M. Smith
and the R Core Team

This manual is for R, version 3.1.2 (2014-10-31).

Copyright © 1990 W. N. Venables

Copyright © 1992 W. N. Venables & D. M. Smith

Copyright © 1997 R. Gentleman & R. Ihaka

Copyright © 1997, 1998 M. Maechler

Copyright © 1999–2014 R Core Team

Permission is granted to make and distribute verbatim copies of this manual provided the copyright notice and this permission notice are preserved on all copies.

Permission is granted to copy and distribute modified versions of this manual under the conditions for verbatim copying, provided that the entire resulting derived work is distributed under the terms of a permission notice identical to this one.

Permission is granted to copy and distribute translations of this manual into another language, under the above conditions for modified versions, except that this permission notice may be stated in a translation approved by the R Core Team.

Table of Contents

Preface	1
1 Introduction and preliminaries	2
1.1 The R environment	2
1.2 Related software and documentation	2
1.3 R and statistics	2
1.4 R and the window system	3
1.5 Using R interactively	3
1.6 An introductory session	4
1.7 Getting help with functions and features	4
1.8 R commands, case sensitivity, etc.	4
1.9 Recall and correction of previous commands	5
1.10 Executing commands from or diverting output to a file	5
1.11 Data permanency and removing objects	5
2 Simple manipulations; numbers and vectors	7
2.1 Vectors and assignment	7
2.2 Vector arithmetic	7
2.3 Generating regular sequences	8
2.4 Logical vectors	9
2.5 Missing values	9
2.6 Character vectors	10
2.7 Index vectors; selecting and modifying subsets of a data set	10
2.8 Other types of objects	11
3 Objects, their modes and attributes	13
3.1 Intrinsic attributes: mode and length	13
3.2 Changing the length of an object	14
3.3 Getting and setting attributes	14
3.4 The class of an object	14
4 Ordered and unordered factors	16
4.1 A specific example	16
4.2 The function <code>tapply()</code> and ragged arrays	16
4.3 Ordered factors	17
5 Arrays and matrices	18
5.1 Arrays	18
5.2 Array indexing. Subsections of an array	18
5.3 Index matrices	19
5.4 The <code>array()</code> function	20
5.4.1 Mixed vector and array arithmetic. The recycling rule	20
5.5 The outer product of two arrays	21
5.6 Generalized transpose of an array	21
5.7 Matrix facilities	22
5.7.1 Matrix multiplication	22

5.7.2	Linear equations and inversion	22
5.7.3	Eigenvalues and eigenvectors	23
5.7.4	Singular value decomposition and determinants	23
5.7.5	Least squares fitting and the QR decomposition	23
5.8	Forming partitioned matrices, <code>cbind()</code> and <code>rbind()</code>	24
5.9	The concatenation function, <code>c()</code> , with arrays	24
5.10	Frequency tables from factors	25
6	Lists and data frames	26
6.1	Lists	26
6.2	Constructing and modifying lists	26
6.2.1	Concatenating lists	27
6.3	Data frames	27
6.3.1	Making data frames	27
6.3.2	<code>attach()</code> and <code>detach()</code>	28
6.3.3	Working with data frames	28
6.3.4	Attaching arbitrary lists	28
6.3.5	Managing the search path	29
7	Reading data from files	30
7.1	The <code>read.table()</code> function	30
7.2	The <code>scan()</code> function	31
7.3	Accessing builtin datasets	31
7.3.1	Loading data from other R packages	31
7.4	Editing data	32
8	Probability distributions	33
8.1	R as a set of statistical tables	33
8.2	Examining the distribution of a set of data	34
8.3	One- and two-sample tests	36
9	Grouping, loops and conditional execution	40
9.1	Grouped expressions	40
9.2	Control statements	40
9.2.1	Conditional execution: <code>if</code> statements	40
9.2.2	Repetitive execution: <code>for</code> loops, <code>repeat</code> and <code>while</code>	40
10	Writing your own functions	42
10.1	Simple examples	42
10.2	Defining new binary operators	43
10.3	Named arguments and defaults	43
10.4	The ‘ <code>...</code> ’ argument	44
10.5	Assignments within functions	44
10.6	More advanced examples	44
10.6.1	Efficiency factors in block designs	44
10.6.2	Dropping all names in a printed array	45
10.6.3	Recursive numerical integration	45
10.7	Scope	46
10.8	Customizing the environment	48
10.9	Classes, generic functions and object orientation	49

11	Statistical models in R	51
11.1	Defining statistical models; formulae	51
11.1.1	Contrasts	53
11.2	Linear models	54
11.3	Generic functions for extracting model information	54
11.4	Analysis of variance and model comparison	55
11.4.1	ANOVA tables	55
11.5	Updating fitted models	55
11.6	Generalized linear models	56
11.6.1	Families	57
11.6.2	The <code>glm()</code> function	57
11.7	Nonlinear least squares and maximum likelihood models	59
11.7.1	Least squares	59
11.7.2	Maximum likelihood	60
11.8	Some non-standard models	61
12	Graphical procedures	63
12.1	High-level plotting commands	63
12.1.1	The <code>plot()</code> function	63
12.1.2	Displaying multivariate data	64
12.1.3	Display graphics	64
12.1.4	Arguments to high-level plotting functions	65
12.2	Low-level plotting commands	66
12.2.1	Mathematical annotation	67
12.2.2	Hershey vector fonts	67
12.3	Interacting with graphics	67
12.4	Using graphics parameters	68
12.4.1	Permanent changes: The <code>par()</code> function	68
12.4.2	Temporary changes: Arguments to graphics functions	69
12.5	Graphics parameters list	69
12.5.1	Graphical elements	70
12.5.2	Axes and tick marks	71
12.5.3	Figure margins	71
12.5.4	Multiple figure environment	73
12.6	Device drivers	74
12.6.1	PostScript diagrams for typeset documents	74
12.6.2	Multiple graphics devices	75
12.7	Dynamic graphics	76
13	Packages	77
13.1	Standard packages	77
13.2	Contributed packages and CRAN	77
13.3	Namespaces	77
14	OS facilities	79
14.1	Files and directories	79
14.2	Filepaths	79
14.3	System commands	80
14.4	Compression and Archives	80
Appendix A	A sample session	82

Appendix B	Invoking R	85
B.1	Invoking R from the command line	85
B.2	Invoking R under Windows	89
B.3	Invoking R under OS X	90
B.4	Scripting with R	90
Appendix C	The command-line editor	92
C.1	Preliminaries	92
C.2	Editing actions	92
C.3	Command-line editor summary	92
Appendix D	Function and variable index	94
Appendix E	Concept index	95
Appendix F	References	96

Preface

This introduction to R is derived from an original set of notes describing the S and S-PLUS environments written in 1990–2 by Bill Venables and David M. Smith when at the University of Adelaide. We have made a number of small changes to reflect differences between the R and S programs, and expanded some of the material.

We would like to extend warm thanks to Bill Venables (and David Smith) for granting permission to distribute this modified version of the notes in this way, and for being a supporter of R from way back.

Comments and corrections are always welcome. Please address email correspondence to R-core@R-project.org.

Suggestions to the reader

Most R novices will start with the introductory session in Appendix A. This should give some familiarity with the style of R sessions and more importantly some instant feedback on what actually happens.

Many users will come to R mainly for its graphical facilities. See [Chapter 12 \[Graphics\]](#), [page 63](#), which can be read at almost any time and need not wait until all the preceding sections have been digested.

1 Introduction and preliminaries

1.1 The R environment

R is an integrated suite of software facilities for data manipulation, calculation and graphical display. Among other things it has

- an effective data handling and storage facility,
- a suite of operators for calculations on arrays, in particular matrices,
- a large, coherent, integrated collection of intermediate tools for data analysis,
- graphical facilities for data analysis and display either directly at the computer or on hard-copy, and
- a well developed, simple and effective programming language (called ‘S’) which includes conditionals, loops, user defined recursive functions and input and output facilities. (Indeed most of the system supplied functions are themselves written in the S language.)

The term “environment” is intended to characterize it as a fully planned and coherent system, rather than an incremental accretion of very specific and inflexible tools, as is frequently the case with other data analysis software.

R is very much a vehicle for newly developing methods of interactive data analysis. It has developed rapidly, and has been extended by a large collection of *packages*. However, most programs written in R are essentially ephemeral, written for a single piece of data analysis.

1.2 Related software and documentation

R can be regarded as an implementation of the S language which was developed at Bell Laboratories by Rick Becker, John Chambers and Allan Wilks, and also forms the basis of the S-PLUS systems.

The evolution of the S language is characterized by four books by John Chambers and coauthors. For R, the basic reference is *The New S Language: A Programming Environment for Data Analysis and Graphics* by Richard A. Becker, John M. Chambers and Allan R. Wilks. The new features of the 1991 release of S are covered in *Statistical Models in S* edited by John M. Chambers and Trevor J. Hastie. The formal methods and classes of the **methods** package are based on those described in *Programming with Data* by John M. Chambers. See [Appendix F \[References\]](#), [page 96](#), for precise references.

There are now a number of books which describe how to use R for data analysis and statistics, and documentation for S/S-PLUS can typically be used with R, keeping the differences between the S implementations in mind. See [Section “What documentation exists for R?”](#) in *The R statistical system FAQ*.

1.3 R and statistics

Our introduction to the R environment did not mention *statistics*, yet many people use R as a statistics system. We prefer to think of it of an environment within which many classical and modern statistical techniques have been implemented. A few of these are built into the base R environment, but many are supplied as *packages*. There are about 25 packages supplied with R (called “standard” and “recommended” packages) and many more are available through the CRAN family of Internet sites (via <http://CRAN.R-project.org>) and elsewhere. More details on packages are given later (see [Chapter 13 \[Packages\]](#), [page 77](#)).

Most classical statistics and much of the latest methodology is available for use with R, but users may need to be prepared to do a little work to find it.

There is an important difference in philosophy between S (and hence R) and the other main statistical systems. In S a statistical analysis is normally done as a series of steps, with intermediate results being stored in objects. Thus whereas SAS and SPSS will give copious output from a regression or discriminant analysis, R will give minimal output and store the results in a fit object for subsequent interrogation by further R functions.

1.4 R and the window system

The most convenient way to use R is at a graphics workstation running a windowing system. This guide is aimed at users who have this facility. In particular we will occasionally refer to the use of R on an X window system although the vast bulk of what is said applies generally to any implementation of the R environment.

Most users will find it necessary to interact directly with the operating system on their computer from time to time. In this guide, we mainly discuss interaction with the operating system on UNIX machines. If you are running R under Windows or OS X you will need to make some small adjustments.

Setting up a workstation to take full advantage of the customizable features of R is a straightforward if somewhat tedious procedure, and will not be considered further here. Users in difficulty should seek local expert help.

1.5 Using R interactively

When you use the R program it issues a prompt when it expects input commands. The default prompt is ‘>’, which on UNIX might be the same as the shell prompt, and so it may appear that nothing is happening. However, as we shall see, it is easy to change to a different R prompt if you wish. We will assume that the UNIX shell prompt is ‘\$’.

In using R under UNIX the suggested procedure for the first occasion is as follows:

1. Create a separate sub-directory, say **work**, to hold data files on which you will use R for this problem. This will be the working directory whenever you use R for this particular problem.

```
$ mkdir work
$ cd work
```

2. Start the R program with the command

```
$ R
```

3. At this point R commands may be issued (see later).
4. To quit the R program the command is

```
> q()
```

At this point you will be asked whether you want to save the data from your R session. On some systems this will bring up a dialog box, and on others you will receive a text prompt to which you can respond **yes**, **no** or **cancel** (a single letter abbreviation will do) to save the data before quitting, quit without saving, or return to the R session. Data which is saved will be available in future R sessions.

Further R sessions are simple.

1. Make **work** the working directory and start the program as before:

```
$ cd work
$ R
```

2. Use the R program, terminating with the **q()** command at the end of the session.

To use R under Windows the procedure to follow is basically the same. Create a folder as the working directory, and set that in the **Start In** field in your R shortcut. Then launch R by double clicking on the icon.

1.6 An introductory session

Readers wishing to get a feel for R at a computer before proceeding are strongly advised to work through the introductory session given in [Appendix A \[A sample session\]](#), page 82.

1.7 Getting help with functions and features

R has an inbuilt help facility similar to the `man` facility of UNIX. To get more information on any specific named function, for example `solve`, the command is

```
> help(solve)
```

An alternative is

```
> ?solve
```

For a feature specified by special characters, the argument must be enclosed in double or single quotes, making it a “character string”: This is also necessary for a few words with syntactic meaning including `if`, `for` and `function`.

```
> help("[")
```

Either form of quote mark may be used to escape the other, as in the string `"It's important"`. Our convention is to use double quote marks for preference.

On most R installations help is available in HTML format by running

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

which will launch a Web browser that allows the help pages to be browsed with hyperlinks. On UNIX, subsequent help requests are sent to the HTML-based help system. The ‘Search Engine and Keywords’ link in the page loaded by `help.start()` is particularly useful as it contains a high-level concept list which searches through available functions. It can be a great way to get your bearings quickly and to understand the breadth of what R has to offer.

The `help.search` command (alternatively `??`) allows searching for help in various ways. For example,

```
> ??solve
```

Try `?help.search` for details and more examples.

The examples on a help topic can normally be run by

```
> example(topic)
```

Windows versions of R have other optional help systems: use

```
> ?help
```

for further details.

1.8 R commands, case sensitivity, etc.

Technically R is an *expression language* with a very simple syntax. It is *case sensitive* as are most UNIX based packages, so `A` and `a` are different symbols and would refer to different variables. The set of symbols which can be used in R names depends on the operating system and country within which R is being run (technically on the *locale* in use). Normally all alphanumeric symbols are allowed¹ (and in some countries this includes accented letters) plus ‘.’ and ‘_’, with the restriction that a name must start with ‘.’ or a letter, and if it starts with ‘.’ the second character must not be a digit. Names are effectively unlimited in length.

Elementary commands consist of either *expressions* or *assignments*. If an expression is given as a command, it is evaluated, printed (unless specifically made invisible), and the value is lost. An assignment also evaluates an expression and passes the value to a variable but the result is not automatically printed.

¹ For portable R code (including that to be used in R packages) only A–Za–z0–9 should be used.

Commands are separated either by a semi-colon (;), or by a newline. Elementary commands can be grouped together into one compound expression by braces ({ and }). *Comments* can be put almost² anywhere, starting with a hashmark (#), everything to the end of the line is a comment.

If a command is not complete at the end of a line, R will give a different prompt, by default

```
+
```

on second and subsequent lines and continue to read input until the command is syntactically complete. This prompt may be changed by the user. We will generally omit the continuation prompt and indicate continuation by simple indenting.

Command lines entered at the console are limited³ to about 4095 bytes (not characters).

1.9 Recall and correction of previous commands

Under many versions of UNIX and on Windows, R provides a mechanism for recalling and re-executing previous commands. The vertical arrow keys on the keyboard can be used to scroll forward and backward through a *command history*. Once a command is located in this way, the cursor can be moved within the command using the horizontal arrow keys, and characters can be removed with the DEL key or added with the other keys. More details are provided later: see [Appendix C \[The command-line editor\]](#), page 92.

The recall and editing capabilities under UNIX are highly customizable. You can find out how to do this by reading the manual entry for the **readline** library.

Alternatively, the Emacs text editor provides more general support mechanisms (via ESS, *Emacs Speaks Statistics*) for working interactively with R. See [Section “R and Emacs” in The R statistical system FAQ](#).

1.10 Executing commands from or diverting output to a file

If commands⁴ are stored in an external file, say `commands.R` in the working directory `work`, they may be executed at any time in an R session with the command

```
> source("commands.R")
```

For Windows **Source** is also available on the **File** menu. The function `sink`,

```
> sink("record.lis")
```

will divert all subsequent output from the console to an external file, `record.lis`. The command

```
> sink()
```

restores it to the console once again.

1.11 Data permanency and removing objects

The entities that R creates and manipulates are known as *objects*. These may be variables, arrays of numbers, character strings, functions, or more general structures built from such components.

During an R session, objects are created and stored by name (we discuss this process in the next session). The R command

```
> objects()
```

(alternatively, `ls()`) can be used to display the names of (most of) the objects which are currently stored within R. The collection of objects currently stored is called the *workspace*.

To remove objects the function `rm` is available:

² **not** inside strings, nor within the argument list of a function definition

³ some of the consoles will not allow you to enter more, and amongst those which do some will silently discard the excess and some will use it as the start of the next line.

⁴ of unlimited length.

```
> rm(x, y, z, ink, junk, temp, foo, bar)
```

All objects created during an R session can be stored permanently in a file for use in future R sessions. At the end of each R session you are given the opportunity to save all the currently available objects. If you indicate that you want to do this, the objects are written to a file called `.RData`⁵ in the current directory, and the command lines used in the session are saved to a file called `.Rhistory`.

When R is started at later time from the same directory it reloads the workspace from this file. At the same time the associated commands history is reloaded.

It is recommended that you should use separate working directories for analyses conducted with R. It is quite common for objects with names `x` and `y` to be created during an analysis. Names like this are often meaningful in the context of a single analysis, but it can be quite hard to decide what they might be when the several analyses have been conducted in the same directory.

⁵ The leading “dot” in this file name makes it *invisible* in normal file listings in UNIX, and in default GUI file listings on OS X and Windows.

2 Simple manipulations; numbers and vectors

2.1 Vectors and assignment

R operates on named *data structures*. The simplest such structure is the numeric *vector*, which is a single entity consisting of an ordered collection of numbers. To set up a vector named `x`, say, consisting of five numbers, namely 10.4, 5.6, 3.1, 6.4 and 21.7, use the R command

```
> x <- c(10.4, 5.6, 3.1, 6.4, 21.7)
```

This is an *assignment* statement using the *function* `c()` which in this context can take an arbitrary number of vector *arguments* and whose value is a vector got by concatenating its arguments end to end.¹

A number occurring by itself in an expression is taken as a vector of length one.

Notice that the assignment operator (`<-`), which consists of the two characters `<` (“less than”) and `-` (“minus”) occurring strictly side-by-side and it ‘points’ to the object receiving the value of the expression. In most contexts the `=` operator can be used as an alternative.

Assignment can also be made using the function `assign()`. An equivalent way of making the same assignment as above is with:

```
> assign("x", c(10.4, 5.6, 3.1, 6.4, 21.7))
```

The usual operator, `<-`, can be thought of as a syntactic short-cut to this.

Assignments can also be made in the other direction, using the obvious change in the assignment operator. So the same assignment could be made using

```
> c(10.4, 5.6, 3.1, 6.4, 21.7) -> x
```

If an expression is used as a complete command, the value is printed *and lost*². So now if we were to use the command

```
> 1/x
```

the reciprocals of the five values would be printed at the terminal (and the value of `x`, of course, unchanged).

The further assignment

```
> y <- c(x, 0, x)
```

would create a vector `y` with 11 entries consisting of two copies of `x` with a zero in the middle place.

2.2 Vector arithmetic

Vectors can be used in arithmetic expressions, in which case the operations are performed element by element. Vectors occurring in the same expression need not all be of the same length. If they are not, the value of the expression is a vector with the same length as the longest vector which occurs in the expression. Shorter vectors in the expression are *recycled* as often as need be (perhaps fractionally) until they match the length of the longest vector. In particular a constant is simply repeated. So with the above assignments the command

```
> v <- 2*x + y + 1
```

generates a new vector `v` of length 11 constructed by adding together, element by element, `2*x` repeated 2.2 times, `y` repeated just once, and `1` repeated 11 times.

The elementary arithmetic operators are the usual `+`, `-`, `*`, `/` and `^` for raising to a power. In addition all of the common arithmetic functions are available. `log`, `exp`, `sin`, `cos`, `tan`, `sqrt`,

¹ With other than vector types of argument, such as `list` mode arguments, the action of `c()` is rather different. See [Section 6.2.1 \[Concatenating lists\]](#), page 27.

² Actually, it is still available as `.Last.value` before any other statements are executed.

and so on, all have their usual meaning. `max` and `min` select the largest and smallest elements of a vector respectively. `range` is a function whose value is a vector of length two, namely `c(min(x), max(x))`. `length(x)` is the number of elements in `x`, `sum(x)` gives the total of the elements in `x`, and `prod(x)` their product.

Two statistical functions are `mean(x)` which calculates the sample mean, which is the same as `sum(x)/length(x)`, and `var(x)` which gives

```
sum((x-mean(x))^2)/(length(x)-1)
```

or sample variance. If the argument to `var()` is an n -by- p matrix the value is a p -by- p sample covariance matrix got by regarding the rows as independent p -variate sample vectors.

`sort(x)` returns a vector of the same size as `x` with the elements arranged in increasing order; however there are other more flexible sorting facilities available (see `order()` or `sort.list()` which produce a permutation to do the sorting).

Note that `max` and `min` select the largest and smallest values in their arguments, even if they are given several vectors. The *parallel* maximum and minimum functions `pmax` and `pmin` return a vector (of length equal to their longest argument) that contains in each element the largest (smallest) element in that position in any of the input vectors.

For most purposes the user will not be concerned if the “numbers” in a numeric vector are integers, reals or even complex. Internally calculations are done as double precision real numbers, or double precision complex numbers if the input data are complex.

To work with complex numbers, supply an explicit complex part. Thus

```
sqrt(-17)
```

will give NaN and a warning, but

```
sqrt(-17+0i)
```

will do the computations as complex numbers.

2.3 Generating regular sequences

R has a number of facilities for generating commonly used sequences of numbers. For example `1:30` is the vector `c(1, 2, ..., 29, 30)`. The colon operator has high priority within an expression, so, for example `2*1:15` is the vector `c(2, 4, ..., 28, 30)`. Put `n <- 10` and compare the sequences `1:n-1` and `1:(n-1)`.

The construction `30:1` may be used to generate a sequence backwards.

The function `seq()` is a more general facility for generating sequences. It has five arguments, only some of which may be specified in any one call. The first two arguments, if given, specify the beginning and end of the sequence, and if these are the only two arguments given the result is the same as the colon operator. That is `seq(2,10)` is the same vector as `2:10`.

Arguments to `seq()`, and to many other R functions, can also be given in named form, in which case the order in which they appear is irrelevant. The first two arguments may be named `from=value` and `to=value`; thus `seq(1,30)`, `seq(from=1, to=30)` and `seq(to=30, from=1)` are all the same as `1:30`. The next two arguments to `seq()` may be named `by=value` and `length=value`, which specify a step size and a length for the sequence respectively. If neither of these is given, the default `by=1` is assumed.

For example

```
> seq(-5, 5, by=.2) -> s3
```

generates in `s3` the vector `c(-5.0, -4.8, -4.6, ..., 4.6, 4.8, 5.0)`. Similarly

```
> s4 <- seq(length=51, from=-5, by=.2)
```

generates the same vector in `s4`.

The fifth argument may be named `along=vector`, which is normally used as the only argument to create the sequence `1, 2, ..., length(vector)`, or the empty sequence if the vector is empty (as it can be).

A related function is `rep()` which can be used for replicating an object in various complicated ways. The simplest form is

```
> s5 <- rep(x, times=5)
```

which will put five copies of `x` end-to-end in `s5`. Another useful version is

```
> s6 <- rep(x, each=5)
```

which repeats each element of `x` five times before moving on to the next.

2.4 Logical vectors

As well as numerical vectors, R allows manipulation of logical quantities. The elements of a logical vector can have the values `TRUE`, `FALSE`, and `NA` (for “not available”, see below). The first two are often abbreviated as `T` and `F`, respectively. Note however that `T` and `F` are just variables which are set to `TRUE` and `FALSE` by default, but are not reserved words and hence can be overwritten by the user. Hence, you should always use `TRUE` and `FALSE`.

Logical vectors are generated by *conditions*. For example

```
> temp <- x > 13
```

sets `temp` as a vector of the same length as `x` with values `FALSE` corresponding to elements of `x` where the condition is *not* met and `TRUE` where it is.

The logical operators are `<`, `<=`, `>`, `>=`, `==` for exact equality and `!=` for inequality. In addition if `c1` and `c2` are logical expressions, then `c1 & c2` is their intersection (“*and*”), `c1 | c2` is their union (“*or*”), and `!c1` is the negation of `c1`.

Logical vectors may be used in ordinary arithmetic, in which case they are *coerced* into numeric vectors, `FALSE` becoming 0 and `TRUE` becoming 1. However there are situations where logical vectors and their coerced numeric counterparts are not equivalent, for example see the next subsection.

2.5 Missing values

In some cases the components of a vector may not be completely known. When an element or value is “not available” or a “missing value” in the statistical sense, a place within a vector may be reserved for it by assigning it the special value `NA`. In general any operation on an `NA` becomes an `NA`. The motivation for this rule is simply that if the specification of an operation is incomplete, the result cannot be known and hence is not available.

The function `is.na(x)` gives a logical vector of the same size as `x` with value `TRUE` if and only if the corresponding element in `x` is `NA`.

```
> z <- c(1:3,NA); ind <- is.na(z)
```

Notice that the logical expression `x == NA` is quite different from `is.na(x)` since `NA` is not really a value but a marker for a quantity that is not available. Thus `x == NA` is a vector of the same length as `x` *all* of whose values are `NA` as the logical expression itself is incomplete and hence undecidable.

Note that there is a second kind of “missing” values which are produced by numerical computation, the so-called *Not a Number*, `NaN`, values. Examples are

```
> 0/0
```

or

both numerical and categorical variables. Many experiments are best described by data frames: the treatments are categorical but the response is numeric. See [Section 6.3 \[Data frames\]](#), page 27.

- *functions* are themselves objects in R which can be stored in the project's workspace. This provides a simple and convenient way to extend R. See [Chapter 10 \[Writing your own functions\]](#), page 42.


```
> winter
```

will print it in data frame form, which is rather like a matrix, whereas

```
> unclass(winter)
```

will print it as an ordinary list. Only in rather special situations do you need to use this facility, but one is when you are learning to come to terms with the idea of class and generic functions.

Generic functions and classes will be discussed further in [Section 10.9 \[Object orientation\]](#), [page 49](#), but only briefly.

as if they were separate vector structures. The result is a structure of the same length as the levels attribute of the factor containing the results. The reader should consult the help document for more details.

Suppose further we needed to calculate the standard errors of the state income means. To do this we need to write an R function to calculate the standard error for any given vector. Since there is a builtin function `var()` to calculate the sample variance, such a function is a very simple one liner, specified by the assignment:

```
> stderr <- function(x) sqrt(var(x)/length(x))
```

(Writing functions will be considered later in [Chapter 10 \[Writing your own functions\]](#), page 42, and in this case was unnecessary as R also has a builtin function `sd()`.) After this assignment, the standard errors are calculated by

```
> incster <- tapply(incomes, statef, stderr)
```

and the values calculated are then

```
> incster
act    nsw  nt    qld    sa tas    vic    wa
1.5 4.3102 4.5 4.1061 2.7386 0.5 5.244 2.6575
```

As an exercise you may care to find the usual 95% confidence limits for the state mean incomes. To do this you could use `tapply()` once more with the `length()` function to find the sample sizes, and the `qt()` function to find the percentage points of the appropriate *t*-distributions. (You could also investigate R's facilities for *t*-tests.)

The function `tapply()` can also be used to handle more complicated indexing of a vector by multiple categories. For example, we might wish to split the tax accountants by both state and sex. However in this simple instance (just one factor) what happens can be thought of as follows. The values in the vector are collected into groups corresponding to the distinct entries in the factor. The function is then applied to each of these groups individually. The value is a vector of function results, labelled by the `levels` attribute of the factor.

The combination of a vector and a labelling factor is an example of what is sometimes called a *ragged array*, since the subclass sizes are possibly irregular. When the subclass sizes are all the same the indexing may be done implicitly and much more efficiently, as we see in the next section.

4.3 Ordered factors

The levels of factors are stored in alphabetical order, or in the order they were specified to **factor** if they were specified explicitly.

Sometimes the levels will have a natural ordering that we want to record and want our statistical analysis to make use of. The `ordered()` function creates such ordered factors but is otherwise identical to **factor**. For most purposes the only difference between ordered and unordered factors is that the former are printed showing the ordering of the levels, but the contrasts generated for them in fitting linear models are different.

5.3 Index matrices

As well as an index vector in any subscript position, a matrix may be used with a single *index matrix* in order either to assign a vector of quantities to an irregular collection of elements in the array, or to extract an irregular collection as a vector.

A matrix example makes the process clear. In the case of a doubly indexed array, an index matrix may be given consisting of two columns and as many rows as desired. The entries in the index matrix are the row and column indices for the doubly indexed array. Suppose for example we have a 4 by 5 array **X** and we wish to do the following:

- Extract elements **X**[1,3], **X**[2,2] and **X**[3,1] as a vector structure, and
- Replace these entries in the array **X** by zeroes.

In this case we need a 3 by 2 subscript array, as in the following example.

```
> x <- array(1:20, dim=c(4,5))    # Generate a 4 by 5 array.
> x
      [,1] [,2] [,3] [,4] [,5]
[1,]     1     5     9    13    17
[2,]     2     6    10    14    18
[3,]     3     7    11    15    19
[4,]     4     8    12    16    20
> i <- array(c(1:3,3:1), dim=c(3,2))
> i                                     # i is a 3 by 2 index array.
      [,1] [,2]
[1,]     1     3
[2,]     2     2
[3,]     3     1
> x[i]                                     # Extract those elements
[1] 9 6 3
> x[i] <- 0                               # Replace those elements by zeros.
> x
      [,1] [,2] [,3] [,4] [,5]
[1,]     1     5     0    13    17
[2,]     2     0    10    14    18
[3,]     0     7    11    15    19
[4,]     4     8    12    16    20
>
```

Negative indices are not allowed in index matrices. **NA** and zero values are allowed: rows in the index matrix containing a zero are ignored, and rows containing an **NA** produce an **NA** in the result.

As a less trivial example, suppose we wish to generate an (unreduced) design matrix for a block design defined by factors **blocks** (**b** levels) and **varieties** (**v** levels). Further suppose there are **n** plots in the experiment. We could proceed as follows:

```
> Xb <- matrix(0, n, b)
> Xv <- matrix(0, n, v)
> ib <- cbind(1:n, blocks)
> iv <- cbind(1:n, varieties)
> Xb[ib] <- 1
> Xv[iv] <- 1
> X <- cbind(Xb, Xv)
```

To construct the incidence matrix, **N** say, we could use

```
> N <- crossprod(Xb, Xv)
```


5.7.3 Eigenvalues and eigenvectors

The function `eigen(Sm)` calculates the eigenvalues and eigenvectors of a symmetric matrix `Sm`. The result of this function is a list of two components named `values` and `vectors`. The assignment

```
> ev <- eigen(Sm)
```

will assign this list to `ev`. Then `ev$val` is the vector of eigenvalues of `Sm` and `ev$vec` is the matrix of corresponding eigenvectors. Had we only needed the eigenvalues we could have used the assignment:

```
> evals <- eigen(Sm)$values
```

`evals` now holds the vector of eigenvalues and the second component is discarded. If the expression

```
> eigen(Sm)
```

is used by itself as a command the two components are printed, with their names. For large matrices it is better to avoid computing the eigenvectors if they are not needed by using the expression

```
> evals <- eigen(Sm, only.values = TRUE)$values
```

5.7.4 Singular value decomposition and determinants

The function `svd(M)` takes an arbitrary matrix argument, `M`, and calculates the singular value decomposition of `M`. This consists of a matrix of orthonormal columns `U` with the same column space as `M`, a second matrix of orthonormal columns `V` whose column space is the row space of `M` and a diagonal matrix of positive entries `D` such that $M = U \%*\% D \%*\% t(V)$. `D` is actually returned as a vector of the diagonal elements. The result of `svd(M)` is actually a list of three components named `d`, `u` and `v`, with evident meanings.

If `M` is in fact square, then, it is not hard to see that

```
> absdetM <- prod(svd(M)$d)
```

calculates the absolute value of the determinant of `M`. If this calculation were needed often with a variety of matrices it could be defined as an R function

```
> absdet <- function(M) prod(svd(M)$d)
```

after which we could use `absdet()` as just another R function. As a further trivial but potentially useful example, you might like to consider writing a function, say `tr()`, to calculate the trace of a square matrix. [Hint: You will not need to use an explicit loop. Look again at the `diag()` function.]

R has a builtin function `det` to calculate a determinant, including the sign, and another, `determinant`, to give the sign and modulus (optionally on log scale),

5.7.5 Least squares fitting and the QR decomposition

The function `lsfit()` returns a list giving results of a least squares fitting procedure. An assignment such as

```
> ans <- lsfit(X, y)
```

gives the results of a least squares fit where `y` is the vector of observations and `X` is the design matrix. See the help facility for more details, and also for the follow-up function `ls.diag()` for, among other things, regression diagnostics. Note that a grand mean term is automatically included and need not be included explicitly as a column of `X`. Further note that you almost always will prefer using `lm(.)` (see [Section 11.2 \[Linear models\]](#), page 54) to `lsfit()` for regression modelling.

Another closely related function is `qr()` and its allies. Consider the following assignments

```

> Xplus <- qr(X)
> b <- qr.coef(Xplus, y)
> fit <- qr.fitted(Xplus, y)
> res <- qr.resid(Xplus, y)

```

These compute the orthogonal projection of y onto the range of X in `fit`, the projection onto the orthogonal complement in `res` and the coefficient vector for the projection in `b`, that is, `b` is essentially the result of the MATLAB ‘backslash’ operator.

It is not assumed that X has full column rank. Redundancies will be discovered and removed as they are found.

This alternative is the older, low-level way to perform least squares calculations. Although still useful in some contexts, it would now generally be replaced by the statistical models features, as will be discussed in [Chapter 11 \[Statistical models in R\], page 51](#).

5.8 Forming partitioned matrices, `cbind()` and `rbind()`

As we have already seen informally, matrices can be built up from other vectors and matrices by the functions `cbind()` and `rbind()`. Roughly `cbind()` forms matrices by binding together matrices horizontally, or column-wise, and `rbind()` vertically, or row-wise.

In the assignment

```
> X <- cbind(arg_1, arg_2, arg_3, ...)
```

the arguments to `cbind()` must be either vectors of any length, or matrices with the same column size, that is the same number of rows. The result is a matrix with the concatenated arguments `arg_1`, `arg_2`, ... forming the columns.

If some of the arguments to `cbind()` are vectors they may be shorter than the column size of any matrices present, in which case they are cyclically extended to match the matrix column size (or the length of the longest vector if no matrices are given).

The function `rbind()` does the corresponding operation for rows. In this case any vector argument, possibly cyclically extended, are of course taken as row vectors.

Suppose $X1$ and $X2$ have the same number of rows. To combine these by columns into a matrix X , together with an initial column of 1s we can use

```
> X <- cbind(1, X1, X2)
```

The result of `rbind()` or `cbind()` always has matrix status. Hence `cbind(x)` and `rbind(x)` are possibly the simplest ways explicitly to allow the vector x to be treated as a column or row matrix respectively.

5.9 The concatenation function, `c()`, with arrays

It should be noted that whereas `cbind()` and `rbind()` are concatenation functions that respect `dim` attributes, the basic `c()` function does not, but rather clears numeric objects of all `dim` and `dimnames` attributes. This is occasionally useful in its own right.

The official way to coerce an array back to a simple vector object is to use `as.vector()`

```
> vec <- as.vector(X)
```

However a similar result can be achieved by using `c()` with just one argument, simply for this side-effect:

```
> vec <- c(X)
```

There are slight differences between the two, but ultimately the choice between them is largely a matter of style (with the former being preferable).

5.10 Frequency tables from factors

Recall that a factor defines a partition into groups. Similarly a pair of factors defines a two way cross classification, and so on. The function `table()` allows frequency tables to be calculated from equal length factors. If there are k factor arguments, the result is a k -way array of frequencies.

Suppose, for example, that `statef` is a factor giving the state code for each entry in a data vector. The assignment

```
> statefr <- table(statef)
```

gives in `statefr` a table of frequencies of each state in the sample. The frequencies are ordered and labelled by the `levels` attribute of the factor. This simple case is equivalent to, but more convenient than,

```
> statefr <- tapply(statef, statef, length)
```

Further suppose that `incomef` is a factor giving a suitably defined “income class” for each entry in the data vector, for example with the `cut()` function:

```
> factor(cut(incomes, breaks = 35+10*(0:7))) -> incomef
```

Then to calculate a two-way table of frequencies:

```
> table(incomef, statef)
```

	statef							
incomef	act	nsw	nt	qld	sa	tas	vic	wa
(35,45]	1	1	0	1	0	0	1	0
(45,55]	1	1	1	1	2	0	1	3
(55,65]	0	3	1	3	2	2	2	1
(65,75]	0	1	0	0	0	0	1	0

Extension to higher-way frequency tables is immediate.


```
> attach(any.old.list)
```

Anything that has been attached can be detached by `detach`, by position number or, preferably, by name.

6.3.5 Managing the search path

The function `search` shows the current search path and so is a very useful way to keep track of which data frames and lists (and packages) have been attached and detached. Initially it gives

```
> search()
[1] ".GlobalEnv"    "Autoloads"     "package:base"
```

where `.GlobalEnv` is the workspace.²

After `lentils` is attached we have

```
> search()
[1] ".GlobalEnv"    "lentils"       "Autoloads"     "package:base"
> ls(2)
[1] "u" "v" "w"
```

and as we see `ls` (or `objects`) can be used to examine the contents of any position on the search path.

Finally, we detach the data frame and confirm it has been removed from the search path.

```
> detach("lentils")
> search()
[1] ".GlobalEnv"    "Autoloads"     "package:base"
```

² See the on-line help for `autoload` for the meaning of the second term.

7.4 Editing data

When invoked on a data frame or matrix, `edit` brings up a separate spreadsheet-like environment for editing. This is useful for making small changes once a data set has been read. The command

```
> xnew <- edit(xold)
```

will allow you to edit your data set `xold`, and on completion the changed object is assigned to `xnew`. If you want to alter the original dataset `xold`, the simplest way is to use `fix(xold)`, which is equivalent to `xold <- edit(xold)`.

Use

```
> xnew <- edit(data.frame())
```

to enter new data via the spreadsheet interface.

Warning message:

cannot compute correct p-values with ties in: `ks.test(A, B)`

10 Writing your own functions

As we have seen informally along the way, the R language allows the user to create objects of mode *function*. These are true R functions that are stored in a special internal form and may be used in further expressions and so on. In the process, the language gains enormously in power, convenience and elegance, and learning to write useful functions is one of the main ways to make your use of R comfortable and productive.

It should be emphasized that most of the functions supplied as part of the R system, such as `mean()`, `var()`, `postscript()` and so on, are themselves written in R and thus do not differ materially from user written functions.

A function is defined by an assignment of the form

```
> name <- function(arg_1, arg_2, ...) expression
```

The *expression* is an R expression, (usually a grouped expression), that uses the arguments, *arg_i*, to calculate a value. The value of the expression is the value returned for the function.

A call to the function then usually takes the form *name*(*expr_1*, *expr_2*, ...) and may occur anywhere a function call is legitimate.

10.1 Simple examples

As a first example, consider a function to calculate the two sample *t*-statistic, showing “all the steps”. This is an artificial example, of course, since there are other, simpler ways of achieving the same end.

The function is defined as follows:

```
> twosam <- function(y1, y2) {
  n1 <- length(y1); n2 <- length(y2)
  yb1 <- mean(y1); yb2 <- mean(y2)
  s1 <- var(y1); s2 <- var(y2)
  s <- ((n1-1)*s1 + (n2-1)*s2)/(n1+n2-2)
  tst <- (yb1 - yb2)/sqrt(s*(1/n1 + 1/n2))
  tst
}
```

With this function defined, you could perform two sample *t*-tests using a call such as

```
> tstat <- twosam(data$male, data$female); tstat
```

As a second example, consider a function to emulate directly the MATLAB backslash command, which returns the coefficients of the orthogonal projection of the vector *y* onto the column space of the matrix, *X*. (This is ordinarily called the least squares estimate of the regression coefficients.) This would ordinarily be done with the `qr()` function; however this is sometimes a bit tricky to use directly and it pays to have a simple function such as the following to use it safely.

Thus given a *n* by 1 vector *y* and an *n* by *p* matrix *X* then *X y* is defined as $(X^T X)^- X^T y$, where $(X^T X)^-$ is a generalized inverse of $X^T X$.

```
> bsplash <- function(X, y) {
  X <- qr(X)
  qr.coef(X, y)
}
```

After this object is created it may be used in statements such as

```
> regcoeff <- bsplash(Xmat, yvar)
```

and so on.


```

area <- function(f, a, b, eps = 1.0e-06, lim = 10) {
  fun1 <- function(f, a, b, fa, fb, a0, eps, lim, fun) {
    ## function 'fun1' is only visible inside 'area'
    d <- (a + b)/2
    h <- (b - a)/4
    fd <- f(d)
    a1 <- h * (fa + fd)
    a2 <- h * (fd + fb)
    if(abs(a0 - a1 - a2) < eps || lim == 0)
      return(a1 + a2)
    else {
      return(fun(f, a, d, fa, fd, a1, eps, lim - 1, fun) +
             fun(f, d, b, fd, fb, a2, eps, lim - 1, fun))
    }
  }
  fa <- f(a)
  fb <- f(b)
  a0 <- ((fa + fb) * (b - a))/2
  fun1(f, a, b, fa, fb, a0, eps, lim, fun1)
}

```

10.7 Scope

The discussion in this section is somewhat more technical than in other parts of this document. However, it details one of the major differences between S-PLUS and R.

The symbols which occur in the body of a function can be divided into three classes; formal parameters, local variables and free variables. The formal parameters of a function are those occurring in the argument list of the function. Their values are determined by the process of *binding* the actual function arguments to the formal parameters. Local variables are those whose values are determined by the evaluation of expressions in the body of the functions. Variables which are not formal parameters or local variables are called free variables. Free variables become local variables if they are assigned to. Consider the following function definition.

```

f <- function(x) {
  y <- 2*x
  print(x)
  print(y)
  print(z)
}

```

In this function, **x** is a formal parameter, **y** is a local variable and **z** is a free variable.

In R the free variable bindings are resolved by first looking in the environment in which the function was created. This is called *lexical scope*. First we define a function called **cube**.

```

cube <- function(n) {
  sq <- function() n*n
  n*sq()
}

```

The variable **n** in the function **sq** is not an argument to that function. Therefore it is a free variable and the scoping rules must be used to ascertain the value that is to be associated with it. Under static scope (S-PLUS) the value is that associated with a global variable named **n**. Under lexical scope (R) it is the parameter to the function **cube** since that is the active binding for the variable **n** at the time the function **sq** was defined. The difference between evaluation

in R and evaluation in S-PLUS is that S-PLUS looks for a global variable called `n` while R first looks for a variable called `n` in the environment created when `cube` was invoked.

```
## first evaluation in S
S> cube(2)
Error in sq(): Object "n" not found
Dumped
S> n <- 3
S> cube(2)
[1] 18
## then the same function evaluated in R
R> cube(2)
[1] 8
```

Lexical scope can also be used to give functions *mutable state*. In the following example we show how R can be used to mimic a bank account. A functioning bank account needs to have a balance or total, a function for making withdrawals, a function for making deposits and a function for stating the current balance. We achieve this by creating the three functions within `account` and then returning a list containing them. When `account` is invoked it takes a numerical argument `total` and returns a list containing the three functions. Because these functions are defined in an environment which contains `total`, they will have access to its value.

The special assignment operator, `<<-`, is used to change the value associated with `total`. This operator looks back in enclosing environments for an environment that contains the symbol `total` and when it finds such an environment it replaces the value, in that environment, with the value of right hand side. If the global or top-level environment is reached without finding the symbol `total` then that variable is created and assigned to there. For most users `<<-` creates a global variable and assigns the value of the right hand side to it². Only when `<<-` has been used in a function that was returned as the value of another function will the special behavior described here occur.

```
open.account <- function(total) {
  list(
    deposit = function(amount) {
      if(amount <= 0)
        stop("Deposits must be positive!\n")
      total <<- total + amount
      cat(amount, "deposited. Your balance is", total, "\n\n")
    },
    withdraw = function(amount) {
      if(amount > total)
        stop("You don't have that much money!\n")
      total <<- total - amount
      cat(amount, "withdrawn. Your balance is", total, "\n\n")
    },
    balance = function() {
      cat("Your balance is", total, "\n\n")
    }
  )
}

ross <- open.account(100)
```

² In some sense this mimics the behavior in S-PLUS since in S-PLUS this operator always creates or assigns to a global variable.

```

robert <- open.account(200)

ross$withdraw(30)
ross$balance()
robert$balance()

ross$deposit(50)
ross$balance()
ross$withdraw(500)

```

10.8 Customizing the environment

Users can customize their environment in several different ways. There is a site initialization file and every directory can have its own special initialization file. Finally, the special functions `.First` and `.Last` can be used.

The location of the site initialization file is taken from the value of the `R_PROFILE` environment variable. If that variable is unset, the file `Rprofile.site` in the R home subdirectory `etc` is used. This file should contain the commands that you want to execute every time R is started under your system. A second, personal, profile file named `.Rprofile`³ can be placed in any directory. If R is invoked in that directory then that file will be sourced. This file gives individual users control over their workspace and allows for different startup procedures in different working directories. If no `.Rprofile` file is found in the startup directory, then R looks for a `.Rprofile` file in the user's home directory and uses that (if it exists). If the environment variable `R_PROFILE_USER` is set, the file it points to is used instead of the `.Rprofile` files.

Any function named `.First()` in either of the two profile files or in the `.RData` image has a special status. It is automatically performed at the beginning of an R session and may be used to initialize the environment. For example, the definition in the example below alters the prompt to `$` and sets up various other useful things that can then be taken for granted in the rest of the session.

Thus, the sequence in which files are executed is, `Rprofile.site`, the user profile, `.RData` and then `.First()`. A definition in later files will mask definitions in earlier files.

```

> .First <- function() {
  options(prompt="$ ", continue="+\t") # $ is the prompt
  options(digits=5, length=999)       # custom numbers and printout
  x11()                               # for graphics
  par(pch = "+")                      # plotting character
  source(file.path(Sys.getenv("HOME"), "R", "mystuff.R"))
                                     # my personal functions
  library(MASS)                      # attach a package
}

```

Similarly a function `.Last()`, if defined, is (normally) executed at the very end of the session. An example is given below.

```

> .Last <- function() {
  graphics.off()                     # a small safety measure.
  cat(paste(date()), "\nAdios\n"))   # Is it time for lunch?
}

```

³ So it is hidden under UNIX.

10.9 Classes, generic functions and object orientation

The class of an object determines how it will be treated by what are known as *generic* functions. Put the other way round, a generic function performs a task or action on its arguments *specific to the class of the argument itself*. If the argument lacks any `class` attribute, or has a class not catered for specifically by the generic function in question, there is always a *default action* provided.

An example makes things clearer. The class mechanism offers the user the facility of designing and writing generic functions for special purposes. Among the other generic functions are `plot()` for displaying objects graphically, `summary()` for summarizing analyses of various types, and `anova()` for comparing statistical models.

The number of generic functions that can treat a class in a specific way can be quite large. For example, the functions that can accommodate in some fashion objects of class `"data.frame"` include

```
[      [[<-    any      as.matrix
[<-    mean     plot     summary
```

A currently complete list can be got by using the `methods()` function:

```
> methods(class="data.frame")
```

Conversely the number of classes a generic function can handle can also be quite large. For example the `plot()` function has a default method and variants for objects of classes `"data.frame"`, `"density"`, `"factor"`, and more. A complete list can be got again by using the `methods()` function:

```
> methods(plot)
```

For many generic functions the function body is quite short, for example

```
> coef
function (object, ...)
  UseMethod("coef")
```

The presence of `UseMethod` indicates this is a generic function. To see what methods are available we can use `methods()`

```
> methods(coef)
[1] coef.aov*           coef.Arima*          coef.default*        coef.listof*
[5] coef.nls*           coef.summary.nls*
```

Non-visible functions are asterisked

In this example there are six methods, none of which can be seen by typing its name. We can read these by either of

```
> getAnywhere("coef.aov")
A single object matching 'coef.aov' was found
It was found in the following places
  registered S3 method for coef from namespace stats
  namespace:stats
with value

function (object, ...)
{
  z <- object$coef
  z[!is.na(z)]
}
```

```
> getS3method("coef", "aov")
function (object, ...)
{
  z <- object$coef
  z[!is.na(z)]
}
```

The reader is referred to the *R Language Definition* for a more complete discussion of this mechanism.

11 Statistical models in R

This section presumes the reader has some familiarity with statistical methodology, in particular with regression analysis and the analysis of variance. Later we make some rather more ambitious presumptions, namely that something is known about generalized linear models and nonlinear regression.

The requirements for fitting statistical models are sufficiently well defined to make it possible to construct general tools that apply in a broad spectrum of problems.

R provides an interlocking suite of facilities that make fitting statistical models very simple. As we mention in the introduction, the basic output is minimal, and one needs to ask for the details by calling extractor functions.

11.1 Defining statistical models; formulae

The template for a statistical model is a linear regression model with independent, homoscedastic errors

$$y_i = \sum_{j=0}^p \beta_j x_{ij} + e_i, \quad e_i \sim \text{NID}(0, \sigma^2), \quad i = 1, \dots, n$$

In matrix terms this would be written

$$y = X\beta + e$$

where the y is the response vector, X is the *model matrix* or *design matrix* and has columns x_0, x_1, \dots, x_p , the determining variables. Very often x_0 will be a column of ones defining an *intercept* term.

Examples

Before giving a formal specification, a few examples may usefully set the picture.

Suppose y , x , x_0 , x_1 , x_2 , \dots are numeric variables, X is a matrix and A , B , C , \dots are factors. The following formulae on the left side below specify statistical models as described on the right.

$y \sim x$

$y \sim 1 + x$ Both imply the same simple linear regression model of y on x . The first has an implicit intercept term, and the second an explicit one.

$y \sim 0 + x$

$y \sim -1 + x$

$y \sim x - 1$ Simple linear regression of y on x through the origin (that is, without an intercept term).

$\log(y) \sim x_1 + x_2$

Multiple regression of the transformed variable, $\log(y)$, on x_1 and x_2 (with an implicit intercept term).

$y \sim \text{poly}(x, 2)$

$y \sim 1 + x + I(x^2)$

Polynomial regression of y on x of degree 2. The first form uses orthogonal polynomials, and the second uses explicit powers, as basis.

$y \sim X + \text{poly}(x, 2)$

Multiple regression y with model matrix consisting of the matrix X as well as polynomial terms in x to degree 2.

$y \sim A$	Single classification analysis of variance model of y , with classes determined by A .
$y \sim A + x$	Single classification analysis of covariance model of y , with classes determined by A , and with covariate x .
$y \sim A*B$	
$y \sim A + B + A:B$	
$y \sim B \%in\% A$	
$y \sim A/B$	Two factor non-additive model of y on A and B . The first two specify the same crossed classification and the second two specify the same nested classification. In abstract terms all four specify the same model subspace.
$y \sim (A + B + C)^2$	
$y \sim A*B*C - A:B:C$	Three factor experiment but with a model containing main effects and two factor interactions only. Both formulae specify the same model.
$y \sim A * x$	
$y \sim A/x$	
$y \sim A/(1 + x) - 1$	Separate simple linear regression models of y on x within the levels of A , with different codings. The last form produces explicit estimates of as many different intercepts and slopes as there are levels in A .
$y \sim A*B + \text{Error}(C)$	An experiment with two treatment factors, A and B , and error strata determined by factor C . For example a split plot experiment, with whole plots (and hence also subplots), determined by factor C .

The operator \sim is used to define a *model formula* in R. The form, for an ordinary linear model, is

$response \sim op_1 term_1 op_2 term_2 op_3 term_3 \dots$

where

response is a vector or matrix, (or expression evaluating to a vector or matrix) defining the response variable(s).

op_i is an operator, either $+$ or $-$, implying the inclusion or exclusion of a term in the model, (the first is optional).

term_i is either

- a vector or matrix expression, or 1,
- a factor, or
- a *formula expression* consisting of factors, vectors or matrices connected by *formula operators*.

In all cases each term defines a collection of columns either to be added to or removed from the model matrix. A 1 stands for an intercept column and is by default included in the model matrix unless explicitly removed.

The *formula operators* are similar in effect to the Wilkinson and Rogers notation used by such programs as Glim and Genstat. One inevitable change is that the operator $'.'$ becomes $':'$ since the period is a valid name character in R.

The notation is summarized below (based on Chambers & Hastie, 1992, p.29):

$Y \sim M$ Y is modeled as M .

$M_1 + M_2$ Include M_1 and M_2 .

$M_1 - M_2$	Include M_1 leaving out terms of M_2 .
$M_1 : M_2$	The tensor product of M_1 and M_2 . If both terms are factors, then the “subclasses” factor.
$M_1 \%in\% M_2$	Similar to $M_1:M_2$, but with a different coding.
$M_1 * M_2$	$M_1 + M_2 + M_1:M_2$.
M_1 / M_2	$M_1 + M_2 \%in\% M_1$.
M^n	All terms in M together with “interactions” up to order n
$I(M)$	Insulate M . Inside M all operators have their normal arithmetic meaning, and that term appears in the model matrix.

Note that inside the parentheses that usually enclose function arguments all operators have their normal arithmetic meaning. The function $I()$ is an identity function used to allow terms in model formulae to be defined using arithmetic operators.

Note particularly that the model formulae specify the *columns of the model matrix*, the specification of the parameters being implicit. This is not the case in other contexts, for example in specifying nonlinear models.

11.1.1 Contrasts

We need at least some idea how the model formulae specify the columns of the model matrix. This is easy if we have continuous variables, as each provides one column of the model matrix (and the intercept will provide a column of ones if included in the model).

What about a k -level factor A ? The answer differs for unordered and ordered factors. For *unordered* factors $k - 1$ columns are generated for the indicators of the second, \dots , k th levels of the factor. (Thus the implicit parameterization is to contrast the response at each level with that at the first.) For *ordered* factors the $k - 1$ columns are the orthogonal polynomials on $1, \dots, k$, omitting the constant term.

Although the answer is already complicated, it is not the whole story. First, if the intercept is omitted in a model that contains a factor term, the first such term is encoded into k columns giving the indicators for all the levels. Second, the whole behavior can be changed by the `options` setting for `contrasts`. The default setting in R is

```
options(contrasts = c("contr.treatment", "contr.poly"))
```

The main reason for mentioning this is that R and S have different defaults for unordered factors, S using Helmert contrasts. So if you need to compare your results to those of a textbook or paper which used S-PLUS, you will need to set

```
options(contrasts = c("contr.helmert", "contr.poly"))
```

This is a deliberate difference, as treatment contrasts (R’s default) are thought easier for newcomers to interpret.

We have still not finished, as the contrast scheme to be used can be set for each term in the model using the functions `contrasts` and `C`.

We have not yet considered interaction terms: these generate the products of the columns introduced for their component terms.

Although the details are complicated, model formulae in R will normally generate the models that an expert statistician would expect, provided that marginality is preserved. Fitting, for example, a model with an interaction but not the corresponding main effects will in general lead to surprising results, and is for experts only.

11.2 Linear models

The basic function for fitting ordinary multiple models is `lm()`, and a streamlined version of the call is as follows:

```
> fitted.model <- lm(formula, data = data.frame)
```

For example

```
> fm2 <- lm(y ~ x1 + x2, data = production)
```

would fit a multiple regression model of y on $x1$ and $x2$ (with implicit intercept term).

The important (but technically optional) parameter `data = production` specifies that any variables needed to construct the model should come first from the `production data frame`. *This is the case regardless of whether data frame `production` has been attached on the search path or not.*

11.3 Generic functions for extracting model information

The value of `lm()` is a fitted model object; technically a list of results of class "lm". Information about the fitted model can then be displayed, extracted, plotted and so on by using generic functions that orient themselves to objects of class "lm". These include

<code>add1</code>	<code>deviance</code>	<code>formula</code>	<code>predict</code>	<code>step</code>
<code>alias</code>	<code>drop1</code>	<code>kappa</code>	<code>print</code>	<code>summary</code>
<code>anova</code>	<code>effects</code>	<code>labels</code>	<code>proj</code>	<code>vcov</code>
<code>coef</code>	<code>family</code>	<code>plot</code>	<code>residuals</code>	

A brief description of the most commonly used ones is given below.

`anova(object_1, object_2)`

Compare a submodel with an outer model and produce an analysis of variance table.

`coef(object)`

Extract the regression coefficient (matrix).

Long form: `coefficients(object)`.

`deviance(object)`

Residual sum of squares, weighted if appropriate.

`formula(object)`

Extract the model formula.

`plot(object)`

Produce four plots, showing residuals, fitted values and some diagnostics.

`predict(object, newdata=data.frame)`

The data frame supplied must have variables specified with the same labels as the original. The value is a vector or matrix of predicted values corresponding to the determining variable values in `data.frame`.

`print(object)`

Print a concise version of the object. Most often used implicitly.

`residuals(object)`

Extract the (matrix of) residuals, weighted as appropriate.

Short form: `resid(object)`.

`step(object)`

Select a suitable model by adding or dropping terms and preserving hierarchies. The model with the smallest value of AIC (Akaike's An Information Criterion) discovered in the stepwise search is returned.

`summary(object)`

Print a comprehensive summary of the results of the regression analysis.

`vcov(object)`

Returns the variance-covariance matrix of the main parameters of a fitted model object.

11.4 Analysis of variance and model comparison

The model fitting function `aov(formula, data=data.frame)` operates at the simplest level in a very similar way to the function `lm()`, and most of the generic functions listed in the table in [Section 11.3 \[Generic functions for extracting model information\]](#), [page 54](#) apply.

It should be noted that in addition `aov()` allows an analysis of models with multiple error strata such as split plot experiments, or balanced incomplete block designs with recovery of inter-block information. The model formula

```
response ~ mean.formula + Error(strata.formula)
```

specifies a multi-stratum experiment with error strata defined by the *strata.formula*. In the simplest case, *strata.formula* is simply a factor, when it defines a two strata experiment, namely between and within the levels of the factor.

For example, with all determining variables factors, a model formula such as that in:

```
> fm <- aov(yield ~ v + n*p*k + Error(farms/blocks), data=farm.data)
```

would typically be used to describe an experiment with mean model $v + n \cdot p \cdot k$ and three error strata, namely “between farms”, “within farms, between blocks” and “within blocks”.

11.4.1 ANOVA tables

Note also that the analysis of variance table (or tables) are for a sequence of fitted models. The sums of squares shown are the decrease in the residual sums of squares resulting from an inclusion of *that term* in the model at *that place* in the sequence. Hence only for orthogonal experiments will the order of inclusion be inconsequential.

For multistratum experiments the procedure is first to project the response onto the error strata, again in sequence, and to fit the mean model to each projection. For further details, see Chambers & Hastie (1992).

A more flexible alternative to the default full ANOVA table is to compare two or more models directly using the `anova()` function.

```
> anova(fitted.model.1, fitted.model.2, ...)
```

The display is then an ANOVA table showing the differences between the fitted models when fitted in sequence. The fitted models being compared would usually be an hierarchical sequence, of course. This does not give different information to the default, but rather makes it easier to comprehend and control.

11.5 Updating fitted models

The `update()` function is largely a convenience function that allows a model to be fitted that differs from one previously fitted usually by just a few additional or removed terms. Its form is

```
> new.model <- update(old.model, new.formula)
```

In the *new.formula* the special name consisting of a period, ‘.’, only, can be used to stand for “the corresponding part of the old model formula”. For example,

```
> fm05 <- lm(y ~ x1 + x2 + x3 + x4 + x5, data = production)
> fm6 <- update(fm05, . ~ . + x6)
> smf6 <- update(fm6, sqrt(.) ~ .)
```

would fit a five variate multiple regression with variables (presumably) from the data frame `production`, fit an additional model including a sixth regressor variable, and fit a variant on the model where the response had a square root transform applied.

Note especially that if the `data=` argument is specified on the original call to the model fitting function, this information is passed on through the fitted model object to `update()` and its allies.

The name ‘.’ can also be used in other contexts, but with slightly different meaning. For example

```
> fmfull <- lm(y ~ . , data = production)
```

would fit a model with response `y` and regressor variables *all other variables in the data frame* `production`.

Other functions for exploring incremental sequences of models are `add1()`, `drop1()` and `step()`. The names of these give a good clue to their purpose, but for full details see the on-line help.

11.6 Generalized linear models

Generalized linear modeling is a development of linear models to accommodate both non-normal response distributions and transformations to linearity in a clean and straightforward way. A generalized linear model may be described in terms of the following sequence of assumptions:

- There is a response, y , of interest and stimulus variables x_1, x_2, \dots , whose values influence the distribution of the response.
- The stimulus variables influence the distribution of y through *a single linear function, only*. This linear function is called the *linear predictor*, and is usually written

$$\eta = \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_p x_p,$$

hence x_i has no influence on the distribution of y if and only if $\beta_i = 0$.

- The distribution of y is of the form

$$f_Y(y; \mu, \varphi) = \exp \left[\frac{A}{\varphi} \{y\lambda(\mu) - \gamma(\lambda(\mu))\} + \tau(y, \varphi) \right]$$

where φ is a *scale parameter* (possibly known), and is constant for all observations, A represents a prior weight, assumed known but possibly varying with the observations, and μ is the mean of y . So it is assumed that the distribution of y is determined by its mean and possibly a scale parameter as well.

- The mean, μ , is a smooth invertible function of the linear predictor:

$$\mu = m(\eta), \quad \eta = m^{-1}(\mu) = \ell(\mu)$$

and this inverse function, $\ell()$, is called the *link function*.

These assumptions are loose enough to encompass a wide class of models useful in statistical practice, but tight enough to allow the development of a unified methodology of estimation and inference, at least approximately. The reader is referred to any of the current reference works on the subject for full details, such as McCullagh & Nelder (1989) or Dobson (1990).

11.6.1 Families

The class of generalized linear models handled by facilities supplied in R includes *gaussian*, *binomial*, *poisson*, *inverse gaussian* and *gamma* response distributions and also *quasi-likelihood* models where the response distribution is not explicitly specified. In the latter case the *variance function* must be specified as a function of the mean, but in other cases this function is implied by the response distribution.

Each response distribution admits a variety of link functions to connect the mean with the linear predictor. Those automatically available are shown in the following table:

Family name	Link functions
binomial	logit, probit, log, cloglog
gaussian	identity, log, inverse
Gamma	identity, inverse, log
inverse.gaussian	1/mu ² , identity, inverse, log
poisson	identity, log, sqrt
quasi	logit, probit, cloglog, identity, inverse, log, 1/mu ² , sqrt

The combination of a response distribution, a link function and various other pieces of information that are needed to carry out the modeling exercise is called the *family* of the generalized linear model.

11.6.2 The glm() function

Since the distribution of the response depends on the stimulus variables through a single linear function *only*, the same mechanism as was used for linear models can still be used to specify the linear part of a generalized model. The family has to be specified in a different way.

The R function to fit a generalized linear model is `glm()` which uses the form

```
> fitted.model <- glm(formula, family=family.generator, data=data.frame)
```

The only new feature is the *family.generator*, which is the instrument by which the family is described. It is the name of a function that generates a list of functions and expressions that together define and control the model and estimation process. Although this may seem a little complicated at first sight, its use is quite simple.

The names of the standard, supplied family generators are given under “Family Name” in the table in [Section 11.6.1 \[Families\], page 57](#). Where there is a choice of links, the name of the link may also be supplied with the family name, in parentheses as a parameter. In the case of the *quasi* family, the variance function may also be specified in this way.

Some examples make the process clear.

The gaussian family

A call such as

```
> fm <- glm(y ~ x1 + x2, family = gaussian, data = sales)
```

achieves the same result as

```
> fm <- lm(y ~ x1+x2, data=sales)
```

but much less efficiently. Note how the gaussian family is not automatically provided with a choice of links, so no parameter is allowed. If a problem requires a gaussian family with a nonstandard link, this can usually be achieved through the *quasi* family, as we shall see later.

The binomial family

Consider a small, artificial example, from Silvey (1970).

Tree models are available in R *via* the user-contributed packages **rpart** and **tree**.

For example, `t()` is the transpose function in R, but users might define their own function named `t`. Namespaces prevent the user's definition from taking precedence, and breaking every function that tries to transpose a matrix.

There are two operators that work with namespaces. The double-colon operator `::` selects definitions from a particular namespace. In the example above, the transpose function will always be available as `base::t`, because it is defined in the `base` package. Only functions that are exported from the package can be retrieved in this way.

The triple-colon operator `:::` may be seen in a few places in R code: it acts like the double-colon operator but also allows access to hidden objects. Users are more likely to use the `getAnywhere()` function, which searches multiple packages.

Packages are often inter-dependent, and loading one may cause others to be automatically loaded. The colon operators described above will also cause automatic loading of the associated package. When packages with namespaces are loaded automatically they are not added to the search list.


```
f <- outer(x, y, function(x, y) cos(y)/(1 + x^2))
```

f is a square matrix, with rows and columns indexed by x and y respectively, of values of the function $\cos(y)/(1 + x^2)$.

```
oldpar <- par(no.readonly = TRUE)
par(pty="s")
```

Save the plotting parameters and set the plotting region to “square”.

```
contour(x, y, f)
contour(x, y, f, nlevels=15, add=TRUE)
```

Make a contour map of f ; add in more lines for more detail.

```
fa <- (f-t(f))/2
```

fa is the “asymmetric part” of f . ($t()$ is transpose).

```
contour(x, y, fa, nlevels=15)
```

Make a contour plot, ...

```
par(oldpar)
```

... and restore the old graphics parameters.

```
image(x, y, f)
image(x, y, fa)
```

Make some high density image plots, (of which you can get hardcopies if you wish), ...

```
objects(); rm(x, y, f, fa)
```

... and clean up before moving on.

R can do complex arithmetic, also.

```
th <- seq(-pi, pi, len=100)
z <- exp(1i*th)
```

$1i$ is used for the complex number i .

```
par(pty="s")
plot(z, type="l")
```

Plotting complex arguments means plot imaginary versus real parts. This should be a circle.

```
w <- rnorm(100) + rnorm(100)*1i
```

Suppose we want to sample points within the unit circle. One method would be to take complex numbers with standard normal real and imaginary parts ...

```
w <- ifelse(Mod(w) > 1, 1/w, w)
```

... and to map any outside the circle onto their reciprocal.

```
plot(w, xlim=c(-1,1), ylim=c(-1,1), pch="+", xlab="x", ylab="y")
lines(z)
```

All points are inside the unit circle, but the distribution is not uniform.

```
w <- sqrt(runif(100))*exp(2*pi*runif(100)*1i)
plot(w, xlim=c(-1,1), ylim=c(-1,1), pch="+", xlab="x", ylab="y")
lines(z)
```

The second method uses the uniform distribution. The points should now look more evenly spaced over the disc.

```
rm(th, w, z)
```

Clean up again.

```
q()
```

Quit the R program. You will be asked if you want to save the R workspace, and for an exploratory session like this, you probably do not want to save it.


```
chem <- scan(n=24)
2.90 3.10 3.40 3.40 3.70 3.70 2.80 2.50 2.40 2.40 2.70 2.20
5.28 3.37 3.03 3.03 28.95 3.77 3.40 2.20 3.50 3.60 3.70 3.70
```

and `stdin()` refers to the script file to allow such traditional usage. If you want to refer to the process's `stdin`, use `"stdin"` as a file connection, e.g. `scan("stdin", ...)`.

Another way to write executable script files (suggested by François Pinard) is to use a *here document* like

```
#!/bin/sh
[environment variables can be set here]
R --slave [other options] <<EOF
```

```
R program goes here...
```

```
EOF
```

but here `stdin()` refers to the program source and `"stdin"` will not be usable.

Short scripts can be passed to `Rscript` on the command-line *via* the `-e` flag. (Empty scripts are not accepted.)

Note that on a Unix-alike the input filename (such as `foo.R`) should not contain spaces nor shell metacharacters.

Appendix C The command-line editor

C.1 Preliminaries

When the GNU **readline** library is available at the time R is configured for compilation under UNIX, an inbuilt command line editor allowing recall, editing and re-submission of prior commands is used. Note that other versions of **readline** exist and may be used by the inbuilt command line editor: this used to happen on OS X.

It can be disabled (useful for usage with ESS¹) using the startup option `--no-readline`.

Windows versions of R have somewhat simpler command-line editing: see ‘**Console**’ under the ‘**Help**’ menu of the GUI, and the file `README.Rterm` for command-line editing under `Rterm.exe`.

When using R with **readline** capabilities, the functions described below are available, as well as others (probably) documented in `man readline` or `info readline` on your system.

Many of these use either Control or Meta characters. Control characters, such as *Control-m*, are obtained by holding the CTRL down while you press the *m* key, and are written as *C-m* below. Meta characters, such as *Meta-b*, are typed by holding down META² and pressing *b*, and written as *M-b* in the following. If your terminal does not have a META key enabled, you can still type Meta characters using two-character sequences starting with *ESC*. Thus, to enter *M-b*, you could type *ESCb*. The *ESC* character sequences are also allowed on terminals with real Meta keys. Note that case is significant for Meta characters.

C.2 Editing actions

The R program keeps a history of the command lines you type, including the erroneous lines, and commands in your history may be recalled, changed if necessary, and re-submitted as new commands. In Emacs-style command-line editing any straight typing you do while in this editing phase causes the characters to be inserted in the command you are editing, displacing any characters to the right of the cursor. In *vi* mode character insertion mode is started by *M-i* or *M-a*, characters are typed and insertion mode is finished by typing a further *ESC*. (The default is Emacs-style, and only that is described here: for *vi* mode see the **readline** documentation.)

Pressing the RET command at any time causes the command to be re-submitted.

Other editing actions are summarized in the following table.

C.3 Command-line editor summary

Command recall and vertical motion

<i>C-p</i>	Go to the previous command (backwards in the history).
<i>C-n</i>	Go to the next command (forwards in the history).
<i>C-r text</i>	Find the last command with the <i>text</i> string in it.

On most terminals, you can also use the up and down arrow keys instead of *C-p* and *C-n*, respectively.

¹ The ‘Emacs Speaks Statistics’ package; see the URL <http://ESS.R-project.org>

² On a PC keyboard this is usually the Alt key, occasionally the ‘Windows’ key. On a Mac keyboard normally no meta key is available.

Horizontal motion of the cursor

<i>C-a</i>	Go to the beginning of the command.
<i>C-e</i>	Go to the end of the line.
<i>M-b</i>	Go back one word.
<i>M-f</i>	Go forward one word.
<i>C-b</i>	Go back one character.
<i>C-f</i>	Go forward one character.

On most terminals, you can also use the left and right arrow keys instead of *C-b* and *C-f*, respectively.

Editing and re-submission

<i>text</i>	Insert <i>text</i> at the cursor.
<i>C-f text</i>	Append <i>text</i> after the cursor.
DEL	Delete the previous character (left of the cursor).
<i>C-d</i>	Delete the character under the cursor.
<i>M-d</i>	Delete the rest of the word under the cursor, and “save” it.
<i>C-k</i>	Delete from cursor to end of command, and “save” it.
<i>C-y</i>	Insert (yank) the last “saved” text here.
<i>C-t</i>	Transpose the character under the cursor with the next.
<i>M-l</i>	Change the rest of the word to lower case.
<i>M-c</i>	Change the rest of the word to upper case.
RET	Re-submit the command to R.

The final RET terminates the command line editing sequence.

The **readline** key bindings can be customized in the usual way *via* a `~/.inputrc` file. These customizations can be conditioned on application R, that is by including a section like

```
$if R
    "\C-xd": "q('no')\n"
$endif
```

Appendix D Function and variable index

!		??	4
!	9	~	
!=	9	~	7
%			
%*%	22	9
%o%	21	40
&		~	
&	9	~	52
&&	40		
*		A	
*	7	abline	66
+		ace	61
+	7	add1	56
-		anova	54, 55
-	7	aov	55
.		aperm	21
.....	55	array	20
.First	48	as.data.frame	27
.Last	48	as.vector	24
/		attach	28
/	7	attr	14
:		attributes	14
:	8	avas	61
::	77	axis	67
:::	77		
<		B	
<	9	boxplot	37
<<-	47	break	41
<=	9	bruto	61
=			
=	9	C	
>		c	7, 10, 24, 27
>	9	C	53
>=	9	cbind	24
?		coef	54
?	4	coefficients	54
		contour	65
		contrasts	53
		coplot	64
		cos	8
		crossprod	19, 22
		cut	25
		D	
		data	31
		data.frame	27
		density	34
		det	23
		detach	28
		determinant	23
		dev.list	75

R

range.....	8
rbind.....	24
read.table.....	30
rep.....	9
repeat.....	41
resid.....	54
residuals.....	54
rlm.....	61
rm.....	6

S

scan.....	31
sd.....	17
search.....	29
seq.....	8
shapiro.test.....	36
sin.....	8
sink.....	5
solve.....	22
sort.....	8
source.....	5
split.....	40
sqrt.....	8
stem.....	34
step.....	54, 56
sum.....	8
summary.....	34, 54
svd.....	23

T

t.....	21
T.....	9
t.test.....	37
table.....	20, 25
tan.....	8
tapply.....	16
text.....	66
title.....	67
tree.....	61
TRUE.....	9

U

unclass.....	14
update.....	55

V

var.....	8, 17
var.test.....	38
vcov.....	55
vector.....	7

W

while.....	41
wilcox.test.....	38
windows.....	74

X

X11.....	74
----------	----

Student's t test 37

T

Tabulation 25

Tree-based models 61

U

Updating fitted models 55

V

Vectors 7

W

Wilcoxon test 38

Workspace 5

Writing functions 42

Appendix F References

- D. M. Bates and D. G. Watts (1988), *Nonlinear Regression Analysis and Its Applications*. John Wiley & Sons, New York.
- Richard A. Becker, John M. Chambers and Allan R. Wilks (1988), *The New S Language*. Chapman & Hall, New York. This book is often called the “*Blue Book*”.
- John M. Chambers and Trevor J. Hastie eds. (1992), *Statistical Models in S*. Chapman & Hall, New York. This is also called the “*White Book*”.
- John M. Chambers (1998) *Programming with Data*. Springer, New York. This is also called the “*Green Book*”.
- A. C. Davison and D. V. Hinkley (1997), *Bootstrap Methods and Their Applications*, Cambridge University Press.
- Annette J. Dobson (1990), *An Introduction to Generalized Linear Models*, Chapman and Hall, London.
- Peter McCullagh and John A. Nelder (1989), *Generalized Linear Models*. Second edition, Chapman and Hall, London.
- John A. Rice (1995), *Mathematical Statistics and Data Analysis*. Second edition. Duxbury Press, Belmont, CA.
- S. D. Silvey (1970), *Statistical Inference*. Penguin, London.