

Computational Statistics

An Introduction to 

Supplement

Günther Sawitzki



CRC Press
Taylor & Francis Group

A CHAPMAN & HALL BOOK

Supplement to
Computational Statistics

An Introduction to 

Günther Sawitzki
StatLab Heidelberg
January 27, 2010

in preparation

<http://sintro.r-forge.r-project.org/>

Introduction

This is a supplement to

G. Sawitzki

Computational Statistics: An Introduction to R

CRC Press, Boca Raton, 2009

It will contain additions, corrections and other supplementary material.

The complete reference appendix of the book is included, with kind permission of CRC Press.

The most recent version of this supplement is on the web site accompanying this book:

<http://sintro.r-forge.r-project.org/>.

Additional material and updates will be available at the web site.

This version: January 27, 2010

Contents

Introduction	v
1 Basic Data Analysis	1
1.5 R Complements	1
1.5.3 Functions	1
Vectorization	1
1.5.6 Search Paths, Frames and Environments	4
2 Regression	9
2.5 Beyond Linear Regression	9
2.5.1 Generalised Linear Models	9
2.6 R Complements	10
2.6.4 Classes and Polymorphic Functions	10
R as a Programming Language and Environment	Suppl. A-11
A.1 Help and Information	Suppl. A-11
A.2 Names and Search Paths	Suppl. A-13
A.3 Administration and Customisation	Suppl. A-15
A.4 Basic Data Types	Suppl. A-16
A.5 Output for Objects	Suppl. A-18
A.6 Object Inspection	Suppl. A-19
A.7 System Inspection	Suppl. A-20
A.8 Complex Data Types	Suppl. A-21
A.9 Accessing Components	Suppl. A-23
A.10 Data Manipulation	Suppl. A-26

A.11 Operators	Suppl. A-29
A.12 Functions	Suppl. A-30
A.13 Debugging and Profiling	Suppl. A-32
A.14 Control Structures	Suppl. A-34
A.15 Input and Output to Data Streams; External Data	Suppl. A-36
A.16 Libraries, Packages	Suppl. A-39
A.17 Mathematical Operators and Functions; Linear Algebra	Suppl. A-41
A.18 Model Descriptions	Suppl. A-42
A.19 Graphic Functions	Suppl. A-44
A.19.1 High-Level Graphics	Suppl. A-44
A.19.2 Low-Level Graphics	Suppl. A-45
A.19.3 Annotations and Legends	Suppl. A-46
A.19.4 Graphic Parameters and Layout	Suppl. A-47
A.20 Elementary Statistical Functions	Suppl. A-48
A.21 Distributions, Random Numbers, Densities...	Suppl. A-49
A.22 Computing on the Language	Suppl. A-52

CHAPTER 1

Basic Data Analysis

1.5 R Complements

1.5.3 Complements: Functions

Vectorization

In R, functions preferably are vectorized. If reasonable, they should accept vectors as parameters, and if appropriate, they should return a vector as result. This is a convenience for calling the function. In some contexts, a function can only be used if it is vectorized. So, for example, `curve()` or `integrate()` only accept a function passed as first argument if it is vectorized.

Unfortunately there is no easy way to check whether a function is vectorized. You must rely on the documentation provided by the author, check the source code, or run some test examples.

If you are providing a function, please document clearly where and how it is vectorized.

You can check whether an argument or any object is a vector using `is.vector()`.

Input _____

```
v <- -1:3
is.vector(v)
```

Output _____

```
[1] TRUE
```

Remember that in R numbers are just vectors of length 1. So `is.vector(7)` will return a `TRUE` value.

Operators in R are functions, and the default operators are vectorized. Most elementary functions are vectorized as well.

Input _____

```
sqrtv<-sqrt(v)
is.vector(sqrtv)
```

[1] TRUE

Input

```
[1]      NaN 0.000000 1.000000 1.414214 1.732051
```

The first example where vectorization usually breaks in your code are logical conditions, because `if` is not vectorized. So the following line takes a vector v , but returns only a vector of length 1.

```
sqrtv <- if (v >= 0) sqrt(v) else 0  
is.vector(sqrtv)
```

Input —

Output
[1] TRUE

Input

[1] 0

If you want to implement a function

$$f(x) = \begin{cases} 0 & x < 0 \\ \sqrt{x} & x \geq 0 \end{cases}$$

the definition

Input –
sqrt0 <- function(x){if (x >=0) sqrt(x) else 0}

does not work as hoped if x is a vector.

In this special case, you can use `ifelse()` which provides a vectorized result.

Exercise 1.1	Vectorization
	Write <code>sqrto()</code> as a vectorized function using <code>ifelse()</code> .

In more general cases, where a vectorized variant is not available, you have to iterate over the components of \mathbf{x} . For performance reasons, the iterators provided with R (Appendix

A.9 Accessing Components (page Suppl. A-23)) are to be preferred over `for`-loop on the indices or other ad-hoc solutions.

Example 1.1: Vectorization with iterators

```
sqrt1 <- function(v) {  
  sapply(v,  
    function(x) {  
      if (x >= 0) sqrt(x) else 0}, simplify=TRUE  
  )  
}
```

Input

To illustrate some technical tricks here, we included the function definition as an anonymous function inline. `sapply()` only relies on the position of the first argument . Argument names do not matter. The elements of `v` are matched to the formal first argument, named `x` in this example.

Of course it is good programming practice to enarmour a general purpose function with guards that check for the appropriate types.

Exercise 1.2	Vectorization
	Enhance <code>sqrt1()</code> above to check that the type of <code>v</code> can be handled correctly. What are the example data types you are using in your test battery?

To facilitate vectorization, a function `Vectorize()` is provided. `Vectorize()` generates an environment (see section Section 1.5.6 (page 4)) and hides your original function there as an object called `FUN`. It then generates a general purpose wrapper to check the arguments and provides vectorization on selected arguments (See Example 1.2 (page 3)).

Example 1.2: Vectorize

<pre>sqrt2 <- Vectorize(sqrt0) sqrt2</pre> <pre>function (x) { args <- lapply(as.list(match.call())[-1L], eval, parent.frame()) names <- if (is.null(names(args))) character(length(args)) else names(args) dovecl <- names %in% vectorize.args do.call("mapply", c(FUN = FUN, args[dovecl], MoreArgs = list(args[!dovecl]), SIMPLIFY = SIMPLIFY, USE.NAMES = USE.NAMES)) }</pre> <p><environment: 0xfd4f20></p> <pre>ls(environment(sqrt2))</pre> <table border="0"> <tr> <td style="vertical-align: top;">[1] "arg.names" "FUN" "FUNV"</td> <td style="vertical-align: top;">[4] "SIMPLIFY" "USE.NAMES" "vectorize.args"</td> </tr> </table> <pre>environment(sqrt2)\$FUN</pre> <pre>function(x){if (x >=0) sqrt(x) else 0}</pre>	[1] "arg.names" "FUN" "FUNV"	[4] "SIMPLIFY" "USE.NAMES" "vectorize.args"	<i>Input</i> _____ <i>Output</i> _____ <i>Input</i> _____ <i>Output</i> _____ <i>Input</i> _____ <i>Output</i> _____ <i>Input</i> _____ <i>Output</i> _____
[1] "arg.names" "FUN" "FUNV"	[4] "SIMPLIFY" "USE.NAMES" "vectorize.args"		

`Vectorize()` is a general purpose tool. If you can provide a special case solution, this may be a chance for optimization. On the long run, general purpose support for vectorization may be enhanced in R. On the other side, more optimizations will be built into the base machine which may make vectorization less critical. But for now, vectorization is a chance for optimization in your code.

1.5.6 Search Paths, Frames and Environments

To evaluate an expression, the formal terms occurring in the expression must be related to actual terms which can be ultimately evaluated. This requires a search process. As the system has evolved, this search process has become rather complex. We try to give a description here, starting with a very simplified picture, and adding details and variations one by one.

In the R documentation, you find several term which are closely related: frames, environments, closures, etc . . . the usage is not always consistent. In particular, *environment* may be used in R documentation where *frame* would be used in older S terminology. In R terminology, environments can be thought of as consisting of two things. A frame, consisting of a set of symbol-value pairs, and an enclosure, a pointer to an enclosing environment. We take the freedom here to define our own usage of these terms.

If you are starting R, several functions and variables are already pre-defined. These come organized in a chain of *environments*. The chain starts with an invisible NULL environment. Next "*package:base*" is the basic environment created upon start of the system. Other environments are populated by loading packages. *search()* gives you a list of the currently active search environments. *searchpaths()* gives you information about the path to the underlying package, if appropriate. Using *ls()*, you can inspect any other environment in this list down to "*package:base*". So *ls("package:base")* gives you a list of the intestines.

For performance reasons, each of these environments is implemented as a data base, called an *environment* in R terminology, and a reference to the predecessor environment. You can think of the data base as a list of names, but actually it contains support for caching and other techniques to improve performance. Moreover this environment does not only contain the name of functions and variables, but it contains name/value pairs.

You start working on the top level. R provides a work space for you, the global environment. Functions and variables you define are added here. This work space environment be accessed as ".*GlobalEnv*" and *ls(".GlobalEnv")* should return a list of your functions and variables. Just calling *ls()* should give the same the same result. *ls.str()* gives you a look at the structure of each entry and its value.

The path can also be modified under program control. For example, a complex data structure like a *data.frame* can be inserted in the search path using *attach()*. After attaching, the components can be found directly. The components are removed from the search path using *detach()*.

Input _____

search()

Output _____

```
[1] ".GlobalEnv"      "tools:RGUI"
[3] "package:stats"   "package:graphics"
[5] "package:grDevices" "package:utils"
[7] "package:datasets" "package:methods"
[9] "Autoloads"        "package:base"
```

Input _____

#ls()

```
expl <- data.frame(x=1:3, y= 11:13) # just an example
```

ls()

Input _____ Output _____

```
[1] "expl"
```

try(x) # component not in search path
try(expl\$x) # component accessible using full name

Output _____

```
[1] 1 2 3
```

attach(expl) Input _____
search() # expl is added after .GlobalEnv

[1] ".GlobalEnv" "expl"
[3] "tools:RGUI" "package:stats"
[5] "package:graphics" "package:grDevices"
[7] "package:utils" "package:datasets"
[9] "package:methods" "Autoloads"
[11] "package:base"

ls() Input _____

[1] "expl" Output _____

try(x) # now in search path

[1] 1 2 3 Output _____

ls("expl") # list objects in the environment attached

[1] "x" "y" Output _____

detach(expl) Input _____
search()

	Output	
[1] ".GlobalEnv"	"tools:RGUI"	
[3] "package:stats"	"package:graphics"	
[5] "package:grDevices"	"package:utils"	
[7] "package:datasets"	"package:methods"	
[9] "Autoloads"	"package:base"	

	Input	
<code>ls()</code>		

	Output	
<code>[1] "expl"</code>		

	Input	
<code>try(x)</code>	# not in search path any more	

Functions are first class objects in R. Functions in R can have formal parameters. They can also have local variables, and functions can be nested in R. As R is an interpreted language, the effective environment can vary. In particular, there is an environment in which a function is defined, and a (usually different) environment in which the function is called. In R, the preference is that variables in a function are evaluated to the values they have when the function is defined. This is called *lexical scoping*. An example is given below.

In more detail, functions have three basic components: a formal argument list, a containing environment, and a body. The combination of these three parts forms what is called the *function closure*. This set defines the lexical scoping. The environment is linking back to the enclosing environment at definition time.

When a variable is requested inside a function, it is first sought in the evaluation environment, then in the enclosure, the enclosure of the enclosure, etc.; once the global environment or the environment of a package is reached, the search continues up the search path to the environment of the base package. If the variable is not found there, the search will proceed next to the empty environment, and will fail.

This construction allows for some optimization. In general, variables are passed by value in R, that is a local copy is generated for each function argument and R functions only operate on the local copy. This causes some time and memory overhead for the copy process. Internally, R uses a lazy evaluation scheme, that is an argument is only evaluated if it is actually needed. Until then, the variable may be treated as a *promise*, defined by an expression to be evaluated, and the environment to be used for evaluation. If the R interpreter recognizes that an argument is unchanged, the copy step may be omitted.

As a special case, environments are never copied, but passed “as is”. So if you have a large data structure, it may be worth considering to hide it as part of the environment, like in the following code fragment which makes use of lexical scoping:

	Input	
--	-------	--

```
definef <- function (x,y, ...) {

  setuphugedata <- function() {
    # some function using x, y, ...
  }
  myhugedata <- setuphugedata()

  return( function () {
    # some function, possibly using myhugedata
    return(myhugedata) # should be some condensed data
  })
}
```

The function `Vectorize()` discussed in Section 1.5.3 (page 1) is an example where this facility is exploited.

After this, `f()` will be a function, accessing `myhugedata` without copy. For a detailed discussion of lexical scoping and more examples, see [3].

At a later stage, a function may be called. This may be from the top level, or from within another environment. When a function is called, a new environment is created, whose enclosure is the environment from the function closure. The run time environment from which a function is called is accessible using `sys.calls()` and `sys.frame()`.

Lexical scoping is the preferred (and default) scoping rule. But expression evaluation is under complete control for the programmer. If you want to use the calling environment as a scope, `sys.frame()` allows to accessing variables and functions by call order, and other rules of scope then R's preferred lexical scoping can be used.

The next detail to add is that on the package level, there is a possibility to fine-tune search paths entries. A package may define a *name space*. Variables and functions are entered into this name space. They may be exported, which will add a reference to the enclosing environment.

With all these possibilities, it is possible that some names are redefined and thus hidden in the search hierarchy. You can however regain hidden definitions by using an explicit reference. So if you have accidentally defined `pi <- 4` and later discovered that the world is not square, you can access the definition of `pi` as given in the base package by using `base::pi`.

CHAPTER 2

Regression**2.5 Beyond Linear Regression***2.5.1 Generalised Linear Models*

We want to proceed to practical work. But at this point we should consider how to overcome the limiting assumptions of linear models. Linear models are among the best-investigated statistical models. Theory and algorithms are far advanced. So it is tempting to try to extend this class of models while still allowing ourselves to use the theoretical and algorithmic know-how.

We have formulated the linear model as

$$\begin{aligned} Y &= m(X) + \varepsilon \\ &\quad Y \text{ with values in } \mathbb{R}^n \\ &\quad X \in \mathbb{R}^{n \times p} \\ &\quad E(\varepsilon) = 0 \\ &\quad \text{with } m(X) = X\beta, \quad \beta \in \mathbb{R}^p. \end{aligned}$$

An important extension is to remove the linearity assumption. As an intermediate step, we do not suppose any longer that m is linear, but only that it can be factored using a linear function. This results in a generalised linear model

$$\begin{aligned} Y &= m(X) + \varepsilon \\ &\quad Y \text{ with values in } \mathbb{R}^n \\ &\quad X \in \mathbb{R}^{n \times p} \\ &\quad E(\varepsilon) = 0 \\ &\quad m(X) = \bar{m}(\eta) \text{ with } \eta = X\beta, \quad \beta \in \mathbb{R}^p. \end{aligned}$$

The next generalisation at hand is to allow for a transformation for Y . Many more generalisations have been discussed. A small number of them have proven tractable. Most important among these is a group of models called generalised linear models (GLM).

An introduction to generalized linear models is [2], and an extensive classical survey is [7].

Generalised linear models have extensive support in R. For most of the functions in R for linear models, there is a corresponding function for generalised linear models. For more information see `help(glm)`.

2.6 R Complements

2.6.4 Classes and Polymorphic Functions

Polymorphism and classes are concepts from object oriented programming. Object oriented programming is a programming style that uses *objects* as basic elements. Objects conceptually consist of *data slots* and *methods*. Encapsulating data and methods as an object is one aspect.

Object oriented programming uses abstract data types, called *classes* which define the data structure and the methods for an object. Classes are arranged in a hierarchy: derived classes inherit the structure of their predecessor, but can add slots and methods. This inheritance is used to generate specific variants. In object oriented programming, variables are instances of a class. The general structure is defined by the class, but the contents of the data slot may be specific to the instance.

Since functions are first class members of R and functions can be stored for example in components of a list, object oriented programming is a style that can be used with R as presented here. Beyond this, R has support for object oriented programming on the language level. `setClass()` allows to define the structure of a class. `new()` is available to create an instance of a class. For details, see chapter 5 of [12].

R as a Programming Language and Environment

R is an interpreted expression language. Expressions are composed of objects and operators.

A.1 Help and Information

Some R functions such as `library()` or `data()` serve a dual purpose. With minimal arguments, they provide help and information. With specific arguments, they give access to certain components.

R Help	
<code>help()</code>	information about an object/a function. <i>Example:</i> <code>help(help)</code>
<code>help.start()</code>	starts browser access to R's online documentation. The reference section includes a search engine to search for keywords, function and data names and text in help page titles.
<code>args()</code>	shows arguments of a function.
<code>example()</code>	executes examples, if available. <i>Example:</i> <code>example(plot)</code>
<code>help.search()</code>	searches for information about an object/a function.
<code>RSiteSearch()</code>	searches for keywords or phrases in the R-help archives or documentation.
<code>apropos()</code>	locates by keyword.
<code>demo()</code>	executes demos for a topic area. <i>Example:</i> <code>demo(graphics)</code> <code>demo()</code> lists all topic areas that provide a demo.

(cont.)→

Suppl. A-12

R Help (cont.)	
library()	gives information about libraries. <i>Example:</i> <code>library()</code> gives a list of all libraries. <code>library(help=<package>)</code> gives information about a package. <i>Example:</i> <code>library(help="stats")</code> gives information about the basic statistics package.
data()	gives information about data sets. <i>Example:</i> <code>data()</code> lists available data sets.
vignette()	lists or views vignette information about a topic. <code>vignette(all = TRUE)</code> lists vignettes from all installed packages. <i>Example:</i> <code>vignette("grid")</code> shows a vignette for the grid graphics.

See also Appendix A.6 “Object Inspection” (page Suppl. A-19) and Appendix A.7 “System Inspection” (page Suppl. A-20).

A.2 Names and Search Paths

Objects are identified by names. By the name objects are searched in a search path, a chain of search environments. The search path in effect can be inspected with `search()`.

R <i>Search Paths</i>	
<code>search()</code>	lists the search areas in effect, beginning with <code>.GlobalEnv</code> down to the base package <code>package:base</code> . <i>Example:</i> <code>search()</code>
<code>searchpaths()</code>	lists the access paths for the search areas in effect. <i>Example:</i> <code>searchpaths()</code>
<code>objects()</code>	lists the objects in a search path. <i>Examples:</i> <code>objects()</code> <code>objects("package:base")</code>
<code>ls()</code>	lists the objects in a search path. <i>Examples:</i> <code>ls()</code> <code>ls("package:base")</code>
<code>ls.str()</code>	lists the objects and their structure in a search path. <i>Examples:</i> <code>ls.str()</code> <code>ls.str("package:base")</code>
<code>find()</code>	locates by keyword. Also finds overlaid entries. <i>Syntax:</i> <code>find(what, mode = "any", numeric = FALSE, simple.words = TRUE)</code>
<code>apropos()</code>	locates by keyword. Also finds overlaid entries. <i>Syntax:</i> <code>apropos(what, where = FALSE, ignore.case = TRUE, mode = "any")</code>

Functions can be nested. This may occur at definition time as well as at execution time. This requires an extension of the search paths. The dynamic identification of objects uses environments to resolve local or global variables in functions.

R <i>Search Paths</i> (cont.)	
<code>environment()</code>	current environments. <i>Example:</i> <code>environment()</code>

(cont.)→

Suppl. A-14

R Search Paths (cont.) (cont.)	
sys.parent()	preceding environments. <i>Example:</i> <code>sys.parent(1)</code>

A.3 Administration and Customisation

<code>objects()</code> <code>ls()</code>	lists the objects in the current search path.
<code>rm()</code>	removes indicated objects. Syntax: <code>rm(<object list>)</code>

R offers a series of possibilities to configure the system so that certain commands are executed upon start or termination. When starting, the files `.Rprofile` and `.RData` are read and executed if available. Details can be system specific. The appropriate information is given by:

```
help(Startup)
```

Various parts of the system keep global information and can be configured by setting options and parameters.

<i>Some System Components with Global State</i>	
basic system	see <code>help(options)</code> .
random numbers	see Appendix A.21 (page Suppl. A-49).
basic graphics	see <code>help(par)</code> .
lattice graphics	see <code>help(lattice.options)</code> .

For information on how to configure memory available for data storage, see:

```
help(Memory)
```

See also Appendix A.7 “System Inspection” (page Suppl. A-20).

A.4 Basic Data Types

<i>Basic R Data Types</i>	
<i>numeric</i>	<i>real</i> or <i>integer</i> . In R: real numbers are always in double precision. Single precision is supported for external call to other languages with .C or .FORTRAN. Functions <i>mode()</i> and <i>typeof()</i> can show the storage modus (single, double ...), depending on the implementations. <i>Examples:</i> <code>1.0</code> <code>2</code> <code>3.14E0</code>
<i>complex</i>	complex, in Cartesian coordinates. <i>Example:</i> <code>1.0+0i</code>
<i>logical</i>	TRUE, FALSE. In R, T and F are predefined variables provided as an alternative. In S-Plus, T and F are basic objects.
<i>character</i>	character strings. Delimiter are alternatively " or '. <i>Example:</i> <code>"T"</code> , <code>'klm'</code>
<i>list</i>	general list structure. List elements can be of different types. <i>Example:</i> <code>list(1:10, "Hello")</code>
<i>function</i>	R function. <i>Example:</i> <code>sin</code>
<i>NULL</i>	special case: empty object. <i>Example:</i> <code>NULL</code>

`is.<type>()` tests for a type, `as.<type>()` converts to a type.

In addition to TRUE and FALSE there are three special values for exceptional situations:

<i>Special Constants</i>	
<i>TRUE</i>	alternative: <i>T</i> . Type: logical.
<i>FALSE</i>	alternative: <i>F</i> . Type: logical.

(cont.)→

<i>Special Constants</i> (cont.)	
<i>NA</i>	“not available”. Type: logical. <i>NA</i> is different from TRUE and FALSE.
<i>NaN</i>	“not a valid numeric value”. Implementation dependent. Should follow the IEEE Standard 754. Type: numeric. <i>Example:</i> 0/0
<i>Inf</i>	infinite. Implementation dependent. Should follow the IEEE Standard 754. Type: numeric. <i>Example:</i> 1/0

<i>Test Functions</i>	
<i>is.na()</i>	returns <i>TRUE</i> if the argument has the value <i>NA</i> or <i>NaN</i> .
<i>na.omit()</i>	returns an object with the cases containing <i>NA</i> removed.
<i>na.fail()</i>	returns its argument if no the case contains <i>NA</i> ; signals an error message otherwise.
<i>is.nan()</i>	returns <i>TRUE</i> if the argument has the value <i>NaN</i> .
<i>is.inf()</i>	returns <i>TRUE</i> if the argument has the value <i>Inf</i> or <i>-Inf</i> .

A.5 Output for Objects

Revised from Appendix ?? (page ??): `str` added.

The object attributes and content can be queried or displayed using output routines. The output routines generally are *polymorphic*, that is they come with variants adapted to the given object type. To list all available methods for a generic function, or all methods for a class, use `methods()`, for example `methods(print)`.

R Inspection	
<code>print()</code>	standard output.
<code>cat()</code>	outputs the objects, concatenating the representations. <code>cat()</code> is useful for producing output in user-defined functions, with minimal formatting.
<code>format()</code>	formats an R object for pretty printing.
<code>structure()</code>	output, optional with attributes.
<code>str()</code>	compact output, optional with attributes.
<code>summary()</code>	standard output as summary, in particular for model fits.
<code>plot()</code>	standard graphic output.

For converting tables to a HTML or LaTeX format, `library(xtable)` [1] is available.

Output of objects to files is discussed in Appendix A.15 “Input and Output to Data Streams” (page Suppl. A-36).

A.6 Object Inspection

Objects have two implicit attributes that can be queried with `mode()` and `length()`. The function `typeof()` gives the (internal) storage modus of an object.

A `class` attribute gives the class of an object.

The following table summarises the most important information possibilities about objects.

<i>Object Inspection</i>	
<code>str()</code>	shows the internal structure of an object in compact form. <i>Syntax:</i> <code>str(<object>)</code>
<code>structure()</code>	shows the internal structure of an object. Attributes for the display can be passed as parameters. <i>Example:</i> <code>structure(1:6, dim = 2:3)</code> <i>Syntax:</i> <code>structure(<object>, ...)</code>
<code>class()</code>	object class. For object classes defined in newer R versions, the class is stored as an attribute. For vintage object classes, the class is determined implicitly by type and other attributes.
<code>mode()</code>	mode (type) of an object.
<code>storage.mode()</code>	storage mode of an object.
<code>typeof()</code>	mode of an object. May be different from the storage mode. Depending on the implementation a numerical variable, for example, can be stored in double precision (the default) or in single precision.
<code>length()</code>	length = number of elements.
<code>attributes()</code>	reads/sets attributes of an object, such as names, dimensions, classes.
<code>names()</code>	names attribute for elements of an object, for example, a vector. <i>Syntax:</i> <code>names(<obj>)</code> gives the <code>names</code> attribute of <code><obj></code> . <code>names(<obj>)<-<charvec></code> sets the <code>names</code> attribute. <i>Example:</i> <code>x<-values</code> <code>names(x)<- <charvec></code>

A.7 System Inspection

The following table summarises the most important information possibilities about the general system environment. When used with an argument, these functions generally serve specific purposes, such as setting parameters and options. When used with an empty argument list, they provide inspection.

<i>System Inspection</i>	
<code>search()</code>	current search path.
<code>ls()</code>	objects in current or selected search path.
<code>methods()</code>	generic methods: Syntax: <code>methods(<fun>)</code> shows specialised functions for <code><fun></code> , <code>methods(class = <c>)</code> the class-specific functions for class <code><c></code> . Examples: <code>methods(plot)</code> <code>methods(class = lm)</code>
<code>data()</code>	accessible data.
<code>library()</code>	accessible packages.
<code>help()</code>	general help system.
<code>options()</code>	global options.
<code>par()</code>	parameter settings for the graphics system.
<code>capabilities()</code>	reports availability of optional features.

The options of the `lattice` systems can be controlled with `trellis.par.set()` resp. `lattice.options()`.

R is anchored in the host operating system. Some variables such as access paths, encoding, etc. are imported from there.

<i>System Environment</i>	
<code>getwd()</code>	gets current working directory.
<code>setwd()</code>	sets current working directory.
<code>dir()</code>	lists files in the current working directory.
<code>system()</code>	calls system functions.

A.8 Complex Data Types

The interpretation of basic types or derived types can be specified by one or more `class` attributes. Polymorphic functions such as `print` or `plot` evaluate this attribute and call a variant for this class if available (see Section 2.6 (page 105)).

For the storage of dates and times, special classes are provided. For more information on these data types see

```
help(DateTimeClasses)
```

and Appendix A.15 (page Suppl. A-36).

R is vector based. Individual constants or values just are vectors with the special length 1. They do not get a special treatment.

<i>Compound Data Types</i>	
Vectors	basic R data types.
Matrices	vectors with two-dimensional layout. See also Appendix A.10 “Data Manipulations” (page Suppl. A-26).
Arrays	vectors with higher-dimensional layout. <code>dim()</code> defines a dimension attribute. <i>Example:</i> <code>x <- runif(100)</code> <code>dim(x) <- c(5, 5, 4)</code> <code>array()</code> generates a new vector with specified dimension structure. <i>Example:</i> <code>z <- array(0, c(4, 3, 2))</code> See also Appendix A.10 “Data Manipulations” (page Suppl. A-26).
Factors	special case for categorical data. <code>factor()</code> converts a numeric vector into a factor. See also Section 2.2. <code>ordered()</code> converts a vector into a factor with ordered levels. This is a shortcut for <code>factor(x, ..., ordered = TRUE)</code> .

(cont.)→

Compound Data Types (cont.)	
	<p>levels() returns the levels of a factor.</p> <p><i>Example:</i> <code>x <- c("a", "b", "a", "c", "a")</code> <code>xf <- factor(x)</code> <code>levels(xf)</code> results in <code>[1] "a" "b" "c"</code></p> <p>tapply() applies a function separately for all levels of factors in a list.</p>
Lists	<p>analogous to vectors, with elements of possibly different types.</p> <p>list() generates a list.</p> <p>Syntax: <code>list(<components>)</code></p> <p><code>[[]]</code> access to components of a list by index.</p> <p><code><list\$component></code> access by names.</p> <p><i>Example:</i> <code>l <- list(name = "xyz", age = 22, fak =</code> <code>"math")</code> <code>>l[[2]]</code> <code>22</code> <code>>l\$age</code> <code>22</code></p>
Data Frames	<p>data frames analogous to arrays resp. lists, with column-wise uniform type and uniform column length.</p> <p>data.frame() analogous to <code>list()</code>, but restrictions have to be satisfied.</p> <p>attach() attaches a database to the current search list. For access to components the component name will be sufficient.</p> <p>detach()</p>

A.9 Accessing Components

The length of vectors is a dynamic attribute. It is extended or shortened as needed. In particular, an implicit “recycling rule” applies: if a vector does not have the length necessary for some operation, it is repeated periodically up to the length required.

Vector components can be accessed by index. The indices can be specified explicitly or in the form of an expression rule.

<i>Accessing Components</i>	
<code>x[<indices>]</code>	indicated components of <i>x</i> . <i>Example:</i> <code>x[1:3]</code>
<code>x[-<indices>]</code>	<i>x</i> omitting indicated components. <i>Example:</i> <code>x[-3]</code> <i>x</i> omitting the 3. component.
<code>x[<condition>]</code>	components of <i>x</i> , for which the <code><condition></code> holds. <i>Example:</i> <code>x[x<0.5]</code>
<code>which()</code>	give the indices of a logical object, allowing for array indices.
<code>subset()</code>	is a polymorphic function and returns subsets of vectors, matrices or data frames by specified conditions.

Vectors (and other objects) can be mapped to higher-dimensional constructs. The layout is described by a additional `dim` attribute. By convention the imbedding goes by column, that is, the first index varies first (FORTRAN convention). Operators and functions can evaluate the dimension attribute.

<i>R Index Access</i>	
<code>dim()</code>	gets or sets dimensions of an object. <i>Example:</i> <code>x <- 1:12; dim(x) <- c(3, 4)</code>
<code>dimnames()</code>	gets or sets names for the dimensions of an object.
<code>nrow()</code>	gives the number of rows = dimension 1.
<code>ncol()</code>	gives the number of columns = dimension 2.
<code>matrix()</code>	generates a matrix with given specifications. <i>Syntax:</i> <code>matrix(data = NA, nrow = 1, ncol = 1, byrow = FALSE, dimnames = NULL)</code> <i>See also</i> Example 1.11 (page 23)

(cont.)→

R Index Access (cont.)	
array()	generates a possibly higher-dimensional matrix. <i>Example:</i> <code>array(x, dim = length(x), dimnames = NULL)</code>

`NCOL()` and `NROW()` are variants treating a vector as a one-column resp. as a one-row matrix.

R Iterators	
apply()	applies a function to the rows or columns of a matrix. <i>Syntax:</i> <code>apply(x, MARGIN, FUNCTION, ...)</code> MARGIN = 1: rows, MARGIN = 2: columns <i>See also</i> Example 1.11 (page 23).
lapply()	applies a function to the elements of a list. <i>Syntax:</i> <code>lapply(X, FUN, ...)</code>
sapply()	applies a function to the elements of a list, of a vector or a matrix. If possible, dimension names are carried over. <i>Syntax:</i> <code>sapply(X, FUN, ..., simplify = TRUE, USE.NAMES = TRUE)</code>
mapply()	applies a function to multiple list or vector arguments. <i>Syntax:</i> <code>mapply(FUN, ..., MoreArgs = NULL, simplify = TRUE, USE.NAMES = TRUE)</code>
Vectorize()	returns a new function that acts as if <code>mapply</code> was called. This can be used as a stepping stone to make a function vectorized. <i>Syntax:</i> <code>Vectorize(FUN, vectorize.args = arg.names, simplify = TRUE, USE.NAMES = TRUE)</code>
tapply()	applies a function to components of an object depending on a list of controlling factors.
by()	object-oriented variant of <code>tapply</code> . <i>Syntax:</i> <code>by(data, INDICES, FUN, ...)</code>
aggregate()	calculates statistics for subsets. <i>Syntax:</i> <code>aggregate(x, ...)</code>

(cont.)→

R <i>Iterators</i> (cont.)	
<code>replicate()</code>	evaluates an expression repeatedly (for example, with generating random numbers for simulation). <i>Syntax:</i> <code>replicate(n, expr, simplify = TRUE)</code>
<code>outer()</code>	generates a matrix with all pair-wise combinations from two vectors, and applies a function to each pair. <i>Syntax:</i> <code>outer(vec1, vec2, FUNCTION, ...)</code>

A.10 Data Manipulation

Array Access	
<code>cbind()</code>	combines by columns.
<code>rbind()</code>	combines by rows.
<code>split()</code>	splits a vector by factors.
<code>table()</code>	generates a table of counts.
<code>prop.table()</code>	expresses table entries as fraction of marginal table, i.e., gives relative counts.
<code>t()</code>	transposes rows and columns. Syntax: <code>t(x)</code>
<code>aperm()</code>	generalised permutation. Syntax: <code>aperm(x, perm)</code> where <code>perm</code> is a permutation of the indices of <code>x</code> .

Transformations	
<code>duplicated()</code>	checks for duplicate or multiple values.
<code>unique()</code>	generates a vector without multiple values.
<code>match()</code>	gives first position of a value in a vector.
<code>pmatch()</code>	partial matching

Character String Transformations	
<code>casefold()</code>	translates characters, in particular from upper - to lowercase or vice versa.
<code>tolower()</code>	translates to lowercase.
<code>toupper()</code>	translates to uppercase.
<code>chartr()</code>	translates characters in a character vector.
<code>substr()</code>	extracts or replaces substrings in a character vector.

(cont.) →

<i>Character String Transformations</i> (cont.)	
substring()	extracts or replaces substrings in a text (respects encoding and other attributes)
paste()	concatenates vectors after converting to character. See also cat() .
strsplit()	splits the elements of a character vector into substrings.
grep()	pattern matching.
gsub()	pattern substitution, by regular patterns.
abbreviate()	abbreviates strings.

<i>Transformations</i>	
table()	generates a table of counts.
expand.grid()	generates a data frame with all combinations of the factors given.
gl()	generates factors by specifying the pattern of their levels.
reshape()	converts between a cross classification table (column per variable) and a long table (variables in rows, with additional indicator column).
merge()	merges data frames. See help(merge) for examples. merge() supports various versions of data base join operations.

<i>Vector Manipulation</i>	
seq()	generates a sequence.
stack()	concatenates multiple vectors from a data frame or list into a single vector and generates a factor indicating the source of each item. <i>Syntax:</i> stack(x, ...)
unstack()	splits a vector by an indicator variable, i.e., reverses the operation of stack() . <i>Syntax:</i> unstack(x, ...)

(cont.)→

Suppl. A-28

<i>Vector Manipulation</i> (cont.)	
<i>split()</i>	splits a vector into the groups defined by a factor. <i>Syntax:</i> <i>split(x, f, drop = FALSE, ...)</i>
<i>unsplit()</i>	combines components to a vector, i.e., reverses <i>split()</i> . <i>Syntax:</i> <i>unsplit(value, f, drop = FALSE)</i>
<i>cut()</i>	converts a numeric to factor. <i>cut()</i> divides the range of a vector into intervals and creates a factor indicating the interval for each value. <i>Syntax:</i> <i>cut(x, ...)</i>

A.11 Operators

Expressions in R can be composed of objects and operators. The following table of operators is ordered by precedence (highest rank on top). See `help(Syntax)`.

<i>Basic R operators</i>	
\$	select component by name. <i>Example:</i> <code>list\$item</code>
[[indexing, access to elements. <i>Example:</i> <code>x[i]</code>
^	exponentiation (right to left). <i>Example:</i> <code>x^3</code>
-	unary minus.
:	sequence generation. <i>Examples:</i> <code>1:5</code> <code>5:1</code>
%<name>%	special operators. Can also be user defined. <i>Examples:</i> <code>"%deg2%"<-function(a, b) a + b^2</code> <code>2 %deg2% 4</code>
* /	multiplication, division.
+ -	addition, subtraction.
< > < = > = ==	comparison operators.
!=	
!	negation.
& &&	and, or . &&, are “shortcut” operators.
<- ->	assignment.

If operands do not have the same length, the shorter operand is repeated cyclically.

Operators of the form `%<name>%` can be defined by the user. The definition follows the rules for function definitions.

Expressions can be written as a sequence with separating semicolons. Expression groups can be combined by enclosing braces `{...}`.

A.12 Functions

Functions are special objects. Functions can return objects as results.

<i>R Function Declarations</i>	
Declarations	<i>function</i> (<i>formal argument list</i>) <i>(expression)</i> <i>Example:</i> <code>fak <- function(n) prod(1:n)</code>
Formal argument	<i>(argument name)</i> <i>(argument name) = (default value)</i>
Formal argument list	list of formal argument, separated by commas. <i>Examples:</i> <code>n, mean = 0, sd = 1</code>
...	variable argument list. Variable argument lists can be propagated to imbedded functions. <i>Example:</i> <code>mean.of.all <- function (...)mean(c(...))</code>
Function result	<i>return</i> <i><value></i> stops function evaluation and returns value. <i><value></i> as last expression in a function declaration: returns value.
Assignments	In general, assignments operate only on local copies of variables. Assignments done within a function are temporary. They are lost after exit from the function. The assignment with <code><<-</code> , however, looks for the target in the complete search chain. It can be used if global and permanent assignments are intended within a function. <i>Syntax:</i> <code><Variable><<-<value></code>

<i>R Function Call</i>	
Function call	<i><name>(<Supplied (actual) argument list>) Example: fak(3)</i>
Supplied argument list	Values are matched by position. Deviating from this, names can be used to control the matching. Initial parts of the names suffice (exception: after a variable argument list, names must be given completely). Function <code>missing()</code> can be used to check, whether a corresponding actual argument is missing for a formal argument. <i>Syntax:</i> <code><list of values> (argument name) = <values></code> <i>Example:</i> <code>rnorm(10, sd = 2)</code>

Arguments for functions are passed by value. This helps consistency, but involves overhead for memory management and copying. If this overhead needs to be avoided, the information provided by `environment()` allows direct access to variables. Techniques to use this are described in [3].

Special case: Functions with names of the form `xxx<-` extend the assignment function. **Example:**

<pre>"inc<-" <-function (x, value) x+value x <- 10 inc(x) <- 3 x</pre>	<i>Input</i> _____
[1] 13	<i>Output</i> _____

In R assignment functions the value argument **must** be called “value”.

A.13 Debugging and Profiling

Revised from Appendix ?? (page ??): changes in “Profiling Support”.

R provides a collection of tools for identification of errors. These are particularly helpful in connection with functions. `browser()` can be used to switch to a browser mode. In this mode, the usual R instructions can be used. Besides this, there is a small number of special instructions. With `debug()`, the browser mode is activated automatically upon entry to a function. The browser mode is marked by a special prompt `Browse[xx]>`.

<code><return></code>	Goes to the next instruction, if the function is under control of <code>debug</code> . Continuous with the expression evaluation if <code>browser</code> has been called directly.
<code>n</code>	Goes to the next instruction (also if <code>browser</code> has been called directly).
<code>cont</code>	Continuous with the expression evaluation.
<code>c</code>	Short for <code>cont</code> . Continues the expression evaluation.
<code>where</code>	Shows call nesting.
<code>Q</code>	Stops execution and jumps back to base state.

Debug Help	
<code>browser()</code>	suspends execution and enters the browser mode. <i>Syntax:</i> <code>browser()</code>
<code>recover()</code>	shows a list of the current call hierarchy. An entry from this list can be chosen for inspection by <code>browser()</code> . With <code>c</code> you leave the <code>browser</code> and return to <code>recover</code> . With <code>0</code> you leave <code>recover()</code> <i>Syntax:</i> <code>recover()</code> <i>Hint:</i> With <code>options(error = recover)</code> , error handling for a function is directed to call <code>browser()</code> automatically in case of an error.
<code>debug()</code>	marks a function for debugger control. On subsequent calls to the function, the debugger is activated and switches to browser mode. <i>Syntax:</i> <code>debug(<function>)</code>
<code>undebug()</code>	cancels debugger control for a function. <i>Syntax:</i> <code>undebug(<function>)</code>
<code>trace()</code>	marks a function for trace control. On subsequent calls to the function, the call is signalled together with its arguments. <i>Syntax:</i> <code>trace(<function>)</code>
<code>untrace()</code>	cancels trace control for a function. <i>Syntax:</i> <code>untrace(<function>)</code>

(cont.)→

Debug Help (cont.)	
traceback()	in case of error inside of a function the current calling stack is stored in a variable <code>.Traceback</code> . <code>traceback()</code> evaluates this variable and displays its content. Syntax: <code>traceback()</code>
try()	Calls a function. Allows for user-defined error handling. Syntax: <code>traceback(<expression>)</code>

To measure execution time in selected code ranges, R provides a “profiling”. This is only available if R has been compiled with the appropriate options. The options installed at compiling can be queried using `capabilities()`. See Appendix A.7 “System Inspection” (page Suppl. A-20).

Profiling Support	
system.time()	returns the execution time of an expression. This function is available in all implementations. Syntax: <code>system.time(<expr>, <gcFirst>)</code>
Rprof()	records active functions periodically. This function is only available if R has been compiled for “profiling”. With <code>memory.profiling = TRUE</code> , in addition to the timing the memory usage is recorded periodically. This option is only available if R has been compiled correspondingly. Syntax: <code>Rprof(filename = "Rprof.out", append = FALSE, interval = 0.02, memory.profiling = FALSE)</code> Use <code>Rprof(NULL)</code> to switch off profiling.
Rprofmem()	records memory requirements on demand. This function is only available if R has been compiled for “memory profiling”. Syntax: <code>Rprofmem(filename = "Rprofmem.out", append = FALSE, threshold = 0)</code> Use <code>Rprofmem(NULL)</code> to switch off profiling.
summaryRprof()	summarises the output of <code>Rprof()</code> and reports the timing by function. Syntax: <code>summaryRprof(filename = "Rprof.out", chunksize = 5000, memory = c("none", "both", "tseries", "stats"), index = 2, diff = TRUE, exclude = NULL)</code> As an alternative, the Perl script <code>R CMD Rprof</code> can be used. See <code>R CMD Rprof -help</code> for usage information.

A.14 Control Structures

R Control Structures	
if	<p>conditional execution.</p> <p>Syntax: <code>if (<log. expression1>) <expression2></code> The logical <code>expression1</code> may return only one logical value. For vector-oriented access use <code>ifelse</code>.</p> <p>Syntax: <code>if ((log. expression1)) <expression2> else <expression3></code></p>
ifelse	<p>element wise conditional execution.</p> <p>Syntax: <code>ifelse(<log. expression1>, <expression2>, <expression3>)</code> evaluates the logical <code>expression1</code> element wise on a vector, and returns <code>expression2</code> if the evaluation gives true, else <code>expression3</code>.</p> <p>Example: <code>trimmedX <- ifelse (abs(X)<2, x, sign(X)*2)</code></p>
switch	<p>evaluates an expression and executes an instruction based on the result.</p> <p>Syntax: <code>switch(<expression1>, ...)</code> <code>expression1</code> must return a numeric value or a character string. ... is an explicit list of alternative actions.</p> <p>Example: <code>centre <- function (x , type) { switch(type, mean = mean(x), median = median(x), trimmed = mean(x, trim = .1))}</code></p>
for	<p>iteration (loop).</p> <p>Syntax: <code>for (<name> in <expression1>) <expression2></code></p>
repeat	<p>iteration. Must be terminated explicitly, for example with <code>break</code>.</p> <p>Syntax: <code>repeat <expression></code></p> <p>Example: <code>pars<-init</code> <code>repeat { res<- get.resid (data, pars)</code> <code>if (converged(res)) break</code> <code>pars<-new.fit (data, pars)}</code></p>

(cont.)→

R Control Structures (cont.)	
while	conditional repetitions. <i>Syntax:</i> <code>while ((log. expression)) <expression></code> <i>Example:</i> <code>pars<-init; res <- get.resid (data, pars)</code> <code>while (!converged(res)) { pars<-</code> <code>new.fit(data, pars)</code> <code>res<- get.resid}</code>
break	terminates the current loop and exits.
next	terminates the current loop cycle and advances to next cycle.

Note: In R loops should be avoided if possible in favour of more efficient language constructs (see [6]).

A.15 Input and Output to Data Streams; External Data

R <i>Input/Output</i>	
<code>write()</code>	writes data to a file. <i>Syntax:</i> <code>write(val, file)</code> <i>Example:</i> <code>write(x, file = "data")</code>
<code>source()</code>	executes the R instruction from the file indicated. <i>Syntax:</i> <code>source("<file name>")</code> <i>Example:</i> <code>source("cmds.R")</code>
<code>Sweave()</code>	executes the R instruction from the file indicated and entangles embedded text. Sweave can be used for automatic report generation. <i>Syntax:</i> <code>Sweave("<file name>", ...)</code>
<code>sink()</code>	redirects output in the file specified. <i>Syntax:</i> <code>sink("<file name>")</code> <i>Example:</i> <code>sink()</code> redirects the output back to the console.
<code>dump()</code>	writes the commands defining an object. The object can be regenerated from this output using <code>source()</code> . <i>Syntax:</i> <code>dump(list, file = "<dumpdata.R>", append = FALSE)</code>

R can access data from local files indicated by a usual file path or from remote files accessed by an URL reference. On most systems, direct access to a clipboard is available as well. More system-specific information is available using `help(connections)`.

To edit or enter data, R provides `edit()`. This is a polymorphic function. For the special case of matrix-like data, `data.entry()` is provided, using a spreadsheet model.

For exchange, the data formats have to be harmonised between all parties. For import from data bases or other systems, several packages are available, for example `library(foreign)` for Stata, SAS, Minitab and SPSS, `library(RODBC)` for SQL. For more information, see the manual “Data Import/Export” [8].

Within R, prepared data are usually provided as *data frames*. If additional objects such as functions or parameters are necessary, they can be made accessible in bundled form as packages. See Appendix A.16 (page Suppl. A-39).

For the exchange from R to R, a special exchange format can be used. Files in this format can be generated with `save()` and conventionally have the name suffix *.Rda*. These files can be loaded again using `load()`.

A general purpose function to load data is `data()`. Depending on the suffix of the input file name, `data()` branches for several special cases. Besides *.Rda* usual suffixes for data input files are *.tab* or *.txt*. The online help function `help(data)` gives additional information.

<i>Data Input/Output for R</i>	
<code>save()</code>	stores data in an external file. Syntax: <code>save(<names of the objects to be stored>, file = <file name>, ...)</code>
<code>save.image()</code>	is a short-cut and stores data of the workspace in an external file.
<code>load()</code>	loads data from an external file. Syntax: <code>load(file = <file name>, ...)</code>
<code>data()</code>	loads data. <code>data()</code> can handle various file formats, if the access paths and filenames follow the R conventions. Syntax: <code>data(..., list = character(0), package = c(.packages(), .Autoloaded), lib.loc = .lib.loc)</code> Example: <code>data(crimes) # loads the data set 'crimes'</code>

For the flexible exchange with other programs in general text-based files are provided. Some conventions can make exchange easier:

- in table form
- only ASCII characters (for example, no umlaut!)
- variables arranged in columns
- columns separated by tabulator stops
- possibly a column header in row 1
- possibly a row label in column 1

For reading the function `read.table()` is provided, and for writing, there is `write.table()`. Besides `read.table()` there are several variants that are adapted to usual data formats. These are documented under `help(read.table)`.

<i>Input and Output of Data for Exchange</i>	
<code>read.table()</code>	reads data tables. Syntax: <code>read.table(file, header = FALSE, sep = "\t", ...)</code> Examples: <code>read.table(<file name>, header = TRUE, sep = "\t")</code> headers in row 1, row labels in column 1 <code>read.table(<file name>, header = TRUE, sep = '\t')</code> now row number, headers in row 1,

(cont.)→

Suppl. A-38

<i>Input and Output of Data for Exchange</i> (cont.)	
<code>write.table()</code>	writes data table. Syntax: <code>write.table(file, header = FALSE, sep = '\t', ...)</code> Examples: <code>write.table(<data frame>, <file name>, header = TRUE, sep = '\t')</code> headers in row 1, row labels in column 1 <code>write.table(<data frame>, <file name>, header = TRUE, sep = '\t')</code> now row number, headers in row 1.
<code>read.csv()</code>	reads comma-separated data tables.
<code>write.csv()</code>	writes comma-separated data tables.
<code>read.csv2()</code>	reads semicolon-separated data tables, using a comma as decimal separator.
<code>write.csv2()</code>	writes semicolon-separated data tables, using a comma as decimal separator.

By default, `read.table()` converts data to `factor` variables if possible. This behaviour can be modified with the argument `as.is` when calling of `read.table()`. This modification is, for example, necessary to read date and time information as for example in the following example from [4]:

```
# date col in all numeric format yyyymmdd
df <- read.table("laketemp.txt", header = TRUE)
as.Date(as.character(df$date), "%Y-%m-%d")
# first two cols in format mm/dd/yy hh:mm:ss
# Note as.is = in read.table to force character
library("chron")
df <- read.table("oxygen.txt", header = TRUE, as.is = 1:2)
chron(df$date, df$time)
```

For sequential reading, `scan()` is provided. Files with data in fixed format (by character columns) can be read with `read.fwf()`.

A.16 Libraries, Packages

External information can be stored in (text) files and packages. In general, additional functions are provided as packages. Packages may be installed as part of the basic installation or installed by the user. Once packages are installed, they are loaded with

`library()`

when needed. Data sets contained in the package are then included in the search path and can be listed using `data()` without arguments:

`data()`

Example:

```
library(nls)
data()
data(Puromycin)
```

If you use R packages, please treat them as you would treat any other scientific source of information. Credit should be given where credit is due, and proper citations should be included. The function `citation()` gives the bibliographic information to use.

<i>Package Utilities</i>	
<code>install.packages()</code>	installs add-on package in <code><lib></code> , downloading it from the archive <i>CRAN</i> or from specified archives. Syntax: <code>install.packages(pkgs, lib, CRAN =getOption("CRAN"), ...)</code> Example: <code>install.packages("mypackage.tgz", repos=NULL)</code> installs package from a local file.
<code>library()</code>	loads an installed add-on package into the current workspace. Syntax: <code>library(package, ...)</code> See also Section 1.5 “Packages” (page 54).
<code>require()</code>	tries to load an add-on package; gives warning on error. Syntax: <code>require(package, ...)</code>
<code>detach()</code>	releases an add-on package and removes it from the search path. Syntax: <code>detach(<name>)</code>
<code>package.manager()</code>	if implemented, interface for management of installed packages. Syntax: <code>package.manager()</code>
<code>package.skeleton()</code>	creates a skeleton for a new package. Syntax: <code>package.skeleton(name = "<anRpackage>", list, ...)</code>
<code>citation()</code>	gives bibliographic information for citing a package. Syntax: <code>citation(<package name>, lib.loc = NULL)</code>

Suppl. A-40

For Unix/Linux/Mac OS X, the main tools are available as commands:

R CMD check <directory> # checks a directory for compliance with the R conventions

R CMD build <directory> # generates an R package

Detailed information for building R packages is in “Writing R Extensions” ([10]).

A.17 Mathematical Operators and Functions; Linear Algebra

For basic arithmetic operators, see `help(Arithmetic)`. For trigonometric functions, information is available using `help(Trig)`. For special mathematical functions, including `beta()`, `factorial()`, `choose()`, see `help(Special)`.

For linear algebra, the most important functions are widely standardised and implemented in C libraries such as BLAS/ATLAS and LAPACK. R makes use of these libraries and provides an interface to the most important functions.

<i>Linear Algebra</i>	
<code>t()</code>	transposes a matrix.
<code>diag()</code>	generates a diagonal matrix.
<code>%^%</code>	matrix multiplication.
<code>rowsum()</code>	gives row sums for a matrix.
<code>colsum()</code>	gives column sums for a matrix.
<code>rowMeans()</code>	gives row means for a matrix.
<code>colMeans()</code>	gives column means for a matrix.
<code>eigen()</code>	computes eigenvalues and eigenvectors of real or complex matrices.
<code>svd()</code>	singular value decomposition of a matrix.
<code>qr()</code>	QR decomposition of a matrix.
<code>determinant()</code>	determinant of a matrix.
<code>solve()</code>	solves linear equations, or computes inverse.

If possible, statistical functions should be used and direct access to the linear algebra functions should be avoided.

<i>Optimisation and Fitting</i>	
<code>optim()</code>	general purpose optimisation.
<code>nlm()</code>	carries out a minimisation of a function using a Newton-type algorithm.
<code>lm()</code>	fits a linear model.
<code>glm()</code>	fits a generalised linear model.
<code>nls()</code>	determines the non-linear (weighted) least-squares estimates of the parameters of a (possibly non-linear) model.
<code>approx()</code>	linear interpolation.
<code>spline()</code>	cubic spline interpolation.

Use the online help functions and search for the keyword `smooth` to find more fitting methods.

A.18 Model Descriptions

Mathematically, linear statistical models can be specified by a design matrix X and written generally as

$$\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \varepsilon,$$

where the matrix \mathbf{X} has to be specified.

R allows us to specify models by giving the rules for how to build the design matrix.

Operator	Syntax	Meaning	Example
\sim	$Y \sim M$	Y depends on M	$Y \sim X$ results in $E(Y) = a + bX$
$+$	$M_1 + M_2$	include M_1 and M_2	$Y \sim X + Z$ $E(Y) = a + bX + cZ$
$-$	$M_1 - M_2$	include M_1 , but exclude M_2	$Y \sim X - 1$ $E(Y) = bX$
:	$M_1 : M_2$	tensor product, that is, all combinations of lev- els of M_1 and M_2	
$\% \text{ in } \%$	$M_1 \% \text{ in } \% M_2$	modified tensor product	$a + b\% \text{ in } \% a$ corre- sponds to $a + a : b$
*	$M_1 * M_2$	“crossed”	$M_1 + M_2$ corre- sponds to $M_1 +$ $M_2 + M_1 : M_2$
/	M_1 / M_2	“nested”: $M_1 + M_2$ $\% \text{ in } \% M_1$	
n	M^n	M with all “interac- tions” up to level n	
$I()$	$I(M)$	interpret M ; terms in M retain their original meaning; the result de- termines the model	$Y \sim (1 + I(X^2))$ corresponds to $E(Y) = a + bX^2$

Table A.34 Wilkinson-Rogers Notation for Linear Models

The model specification is also possible for generalised (not linear) models.

Examples:

$$\begin{aligned} y &\sim 1 + x \\ y &\sim x \end{aligned} \quad \begin{aligned} &\text{corresponds to } y_i = (1 \ x_i)(\beta_1 \ \beta_2)^\top + \varepsilon \\ &\text{short for } y \sim 1 + x \\ &\quad (\text{a constant term is assumed implicitly}) \end{aligned}$$

$y \sim 0 + x$	corresponds to $y_i = x_i \cdot \beta + \varepsilon$
$\log(y) \sim x_1 + x_2$	corresponds to $\log(y_i) = (1 \ x_{i1} \ x_{i2})(\beta_1 \ \beta_2 \ \beta_3)^\top + \varepsilon$ (a constant term is assumed implicitly)
$y \sim A$	one-way analysis of variance with factor A
$y \sim A + x$	covariance analysis with factor A and covariable x
$y \sim A * B$	two-factor crossed layout with factors A and B
$y \sim A/B$	two-factor hierarchical layout with factor A and sub-factor B

For an economic transition between models, for example for model comparison, the function `update()` is available. It updates and (by default) re-fits a model by extracting the call stored in the object, updating the call and evaluating that call, given the new information. In particular, it can be used to re-fit a model to a changed (possibly corrected) data set.

<i>Model Administration</i>	
<code>formula()</code>	extracts a model formula from an object.
<code>terms()</code>	extracts terms of the model formula from an object.
<code>contrasts()</code>	specifies contrasts.
<code>update()</code>	updates and re-fits, or changes a model.
<code>model.matrix()</code>	generates the design matrix for a model.

Example:

```
lm(y ~ poly(x, 4), data = experiment)
```

analyses the data set “experiment” with a linear model for polynomial regression of degree 4.

<i>Standard Analysis</i>	
<code>lm()</code>	linear model. See also Chapter 2.
<code>glm()</code>	generalised linear model.
<code>nls()</code>	non-linear least squares.
<code>nlm()</code>	general non-linear minimisation.
<code>update()</code>	update and re-fit, or change a model.
<code>anova()</code>	analysis of variance.

A.19 Graphic Functions

R provides two graphics systems: The basic graphics system of R implements a model that is oriented at pen and paper drawing. The lattice graphics system is an additional second graphics system that is oriented at a viewport/object model. For information about lattice see `help(lattice)`. For a survey about the functions in lattice see `library(help = lattice)`. Information about the basic graphics system follows here. Additional graphics systems are available as packages.

Graphic functions fall essentially in three groups:

- “high level” functions. These define a new output.
- “low level” functions. These modify an existing output.
- parametrisations. These modify the settings of the graphics system.

A.19.1 High-Level Graphics

High-Level Graphics	
<code>plot()</code>	generic graphic output.
<code>pairs()</code>	pair-wise scatterplots.
<code>coplot()</code>	scatterplots, conditioned on covariates.
<code>qqplot()</code>	<i>QQ</i> Plot.
<code>qqnorm()</code>	Gaussian <i>QQ</i> Plot.
<code>qqline()</code>	adds a line to a Gaussian <i>QQ</i> Plot, passing through the first and third quartile.
<code>hist()</code>	histogram. See also Section 1.3, page 16.
<code>boxplot()</code>	box-and-whisker plot.
<code>dotchart()</code>	draws a Cleveland dot plot.
<code>curve()</code>	evaluates a function or an expression and draws a curve. <i>Example:</i> <code>curve(dnorm, from = -3, to = 3)</code>
<code>image()</code>	colour coded z against x, y .
<code>contour()</code>	contour plot of z against x, y .
<code>persp()</code>	3D surface.
<code>matplot()</code>	plots the columns of one matrix against the columns of another.
<code>mosaicplot()</code>	mosaic displays to visualise (standardised) residuals of a log-linear model for the table.
<code>termplot()</code>	plots regression terms against their predictors, optionally with standard errors and partial residuals added.

Corresponding function names for the `lattice` graphics are in Table 4.5 (page 147).

A.19.2 Low-Level Graphics

Most high-level functions have an argument `add`. If the function is called with `add = FALSE`, it can be used to add elements to an existing plot. Moreover, there are several low-level functions that suppose that there is already a defined plot environment. This is usually set by high-level functions, but may be modified by `par()`: Besides the physical layout, information about the scales, such as range and possible logarithmic transformations, are part of the environment.

Low-Level Plotting	
<code>points()</code>	generic function. Marks points at specified positions. Syntax: <code>points(x, ...)</code>
<code>symbols()</code>	draws symbols at selected points.
<code>text()</code>	adds text labels at selected points.
<code>lines()</code>	generic function. Joins points at specified positions. Syntax: <code>lines(x, ...)</code>
<code>segments()</code>	adds line segments.
<code>abline()</code>	adds a line (in several representations) to a plot. Syntax: <code>abline(a, b, ...)</code>
<code>arrows()</code>	adds a line with arrows to a plot.
<code>polygon()</code>	adds polygon with specified vertices.
<code>rect()</code>	draws a rectangle.
<code>axis()</code>	adds axis.
<code>rug()</code>	adds a rug marking the data points.

Besides this, R has rudimentary possibilities for interaction with graphics.

Interactions	
<code>devAskNewPage()</code>	controls if a console prompt is given before starting a page of output.
<code>locator()</code>	determines the position of mouse clicks. A current graphics display has to be defined before <code>locator()</code> is used. Example: <code>plot(runif(19))</code> <code>locator(n = 3, type = "l")</code>
<code>Sys.sleep()</code>	suspends execution for a time interval. Syntax: <code>Sys.sleep(<seconds>)</code>

(cont.)→

<i>Interactions</i> (cont.)	
<code>getGraphicsEvent()</code>	waits for a keyboard or mouse event. Functions to respond to these events can be specified. This function needs a graphics display that supports graphics events.

For more interactive facilities, see additional packages, in particular:

- rgl implements OpenGL for real-time 3d rendering,
- rggobi interfaces to the ggobi system for higher-dimensional exploration of data.

A.19.3 Annotations and Legends

The high-level functions generally offer the possibilities to add standard annotations by using arguments:

```
main =      main title, above the plot,
sub =       plot caption, below the plot,
xlab =      label for the x axis,
ylab =      label for the y axis.
```

For documentation, see `help(plot.default)`.

High-level functions are complemented by low-level functions.

<i>Low Level Annotation</i>	
<code>title()</code>	adds main title, analogous to high-level argument <code>main</code> . <i>Syntax:</i> <code>title(main = NULL, sub = NULL, xlab = NULL, ylab = NULL, ...)</code>
<code>text()</code>	adds text at specified coordinates. <i>Syntax:</i> <code>text(x, y = NULL, text, ...)</code>
<code>legend()</code>	adds a legend block. <i>Syntax:</i> <code>legend(x, y = NULL, text, ...)</code> <i>Example:</i> <code>plot(runif(100)); legend(locator(1), legend = "You clicked here")</code>
<code>mtext()</code>	adds text to margin. <i>Syntax:</i> <code>mtext(text, side = 3, ...)</code> . The margins are denoted by 1 = bottom, 2 = left, 3 = top, 4 = right).

For annotations, texts some times has to be shortened. Function and variable names can be shortened using `abbreviate()`.

R gives (limited) possibilities for mathematical typesetting. If the text argument is a character string, it is taken directly. If the text argument is an (unevaluated) R expression, R tries to render the expression as usual in a mathematical formula. R expressions can be generated using the functions `expression()` and evaluated with `eval()` or `bquote()`.

Example:

```
text(x, y, expression(paste(bquote("(",
  atop(n, x), ")"),
  .(p)\^{}x, .(q)\^{}n-x\})))
```

`demo(plotmath)` gives several examples for mathematical typesetting in plots.

A.19.4 Graphic Parameters and Layout

<i>Parametrisations</i>	
<code>par()</code>	sets parameters for the basic graphics system. <i>Syntax:</i> see <code>help(par)</code> . <i>Example:</i> <code>par(mfrow = c(m, n))</code> splits the graphic area in <i>m</i> rows and <i>n</i> columns, to be filled row-wise. <code>par(mfcol = c(m, n))</code> fills the area column by column.
<code>lattice.options()</code>	sets parameters for the lattice graphics system. <i>Syntax:</i> see <code>help(lattice.options)</code>
<code>split.screen()</code>	splits the graphic area in parts. <i>Syntax:</i> <code>split.screen(figs, screen, erase = TRUE)</code> . If <i>figs</i> is a pair of two arguments, these will fix the number of rows and columns. If <i>figs</i> is a matrix, each row gives the coordinates of a graphic area in relative coordinates [0...1]. <code>split.screen()</code> can be nested.
<code>screen()</code>	selects graphic area for the next graphical output. <i>Syntax:</i> <code>screen(n = cur.screen, new = TRUE)</code> .
<code>layout()</code>	divides the graphic area. This function is not compatible with other layout functions.

A.20 Elementary Statistical Functions

<i>Statistical Functions</i>	
<i>sum()</i>	sums up components of a vector.
<i>cumsum()</i>	calculates cumulated sums.
<i>prod()</i>	multiplies components of a vector.
<i>cumprod()</i>	calculates cumulated products.
<i>length()</i>	length of an object, for example a vector.
<i>max()</i>	maximum, minimum.
<i>min()</i>	See also <i>pmax</i> , <i>pmin</i> .
<i>range()</i>	minimum and maximum.
<i>cummax()</i>	cumulated maximum, minimum.
<i>cummin()</i>	
<i>quantile()</i>	sample quantile. For theoretical distributions, use <i>qxxxx</i> , for example <i>qnorm</i> .
<i>median()</i>	median.
<i>mean()</i>	mean, including trimmed mean.
<i>var()</i>	variance, variance / covariance matrix.
<i>sort()</i>	sorting.
<i>rev()</i>	reverse sorting.
<i>order()</i>	returns a permutation for sorting.
<i>which.max()</i>	index of the (first) maximum of a numeric vector.
<i>which.min()</i>	index of the (first) minimum of a numeric vector.
<i>rank()</i>	rank in a sample.

A.21 Distributions, Random Numbers, Densities...

The base generator for uniform random numbers is administered by `Random`. Several types of generators are available as base generator. **For serious simulation it is strongly recommended to read the recommendations of Marsaglia et al.** (see `help(.Random.seed)`). All non-uniform random number generators are derived from the current base generator. A survey of most important non-uniform random number generators, their distribution functions and their quantiles is given at the end of this section.

R Random Numbers	
<code>.Random.seed</code>	<p><code>.Random.seed</code> is a global variable that holds the current state of the basic random number generator. This variable can be stored and later be restored with <code>set.seed()</code>.</p> <p>Initially, there is no seed. Use <code>set.seed()</code> to define a seed. If no seed has been defined, a new one is created based on the current clock time when one is required.</p> <p>Random number generators may use variables other than <code>.Random.seed</code> to store their state information. To set a generator to a defined state, always use <code>set.seed()</code>. Never set <code>.Random.seed</code> directly.</p>
<code>set.seed()</code>	<p>initialises the random number generator.</p> <p>Syntax: <code>set.seed(seed, kind = NULL)</code></p>
<code>RNGkind()</code>	<p><code>RNGkind()</code> gives the name of the current base generator.</p> <p><code>RNGkind(<name>)</code> sets a basic random number generator.</p> <p>Syntax: <code>RNGkind()</code> <code>RNGkind(<name>)</code></p> <p>Example: <code>RNGkind("Wichmann-Hill")</code> <code>RNGkind("Marsaglia-Multicarry")</code> <code>RNGkind("Super-Duper")</code></p>
<code>sample()</code>	<p>draws a sample from the values given in vector <code>x</code>, with or without replacement (controlled by the value of <code>replace</code>).</p> <p>Size is by default the length of <code>x</code>.</p> <p>Optionally, <code>prob</code> can be a vector of probabilities for the values of <code>x</code>.</p> <p>Syntax: <code>sample(x, size, replace = FALSE, prob)</code></p> <p>Example: Random permutation: <code>sample(x)</code></p> <p>Biased coin: <code>val<-c("H", "T")</code> <code>prob<-c(0.3, 0.7)</code> <code>sample(val, 10, replace = TRUE, prob)</code></p>

Suppl. A-50

If simulations shall be reproducible, the random number generator must be set to a well-defined initial state for a reproduction. So the initial state needs to be recorded. An example is the following statement sequence to store the current state:

```
save.seed <- .Random.seed
save.kind <- RNGkind()
```

These variables can be stored to a file and read from there when necessary. With

```
set.seed(save.seed, save.kind)
```

the state of the random number generator is then restored.

The individual function names for the common non-uniform generators and distribution functions are combined from a prefix and the short name of the distribution (see the list below). General pattern: if *xxxx* is the short name, then

- rxxxx* generates random numbers
- dxxxx* density or probability
- pxxxx* distribution function
- qxxxx* quantiles

Example:

- `x<-runif(100)` generates 100 random variables with $U(0, 1)$ distribution.
- `qf(0.95, 10, 2)` calculates the 95% quantile of the $F(10, 2)$ distribution.

Distributions	Short Name	Parameter and Default Values
Beta	<code>beta</code>	<code>shape1, shape2, ncp = 0</code>
Binomial	<code>binom</code>	<code>size, prob</code>
Cauchy	<code>cauchy</code>	<code>location = 0, scale = 1</code>
χ^2	<code>chisq</code>	<code>df, ncp = 0</code>
Exponential	<code>exp</code>	<code>rate = 1</code>
F	<code>f</code>	<code>df1, df2 (ncp = 0)</code>
Gamma	<code>gamma</code>	<code>shape, scale = 1</code>
Gauss	<code>norm</code>	<code>mean = 0, sd = 1</code>
Geometric	<code>geom</code>	<code>prob</code>
Hypergeometric	<code>hyper</code>	<code>m, n, k</code>
Lognormal	<code>lnorm</code>	<code>meanlog = 0, sdlog = 1</code>
Logistic	<code>logis</code>	<code>location = 0, scale = 1</code>
Negativ-Binomial	<code>nbinom</code>	<code>size, prob</code>
Poisson	<code>pois</code>	<code>lambda</code>
Student's t	<code>t</code>	<code>df</code>

(cont.)→

<i>Distributions</i>	<i>Short Name</i>	<i>Parameter and Default Values</i>
Tukey Studentised Range	<i>tukey</i>	
Uniform	<i>unif</i>	<i>min = 0, max = 1</i>
Wilcoxon Signed Rank	<i>signrank</i>	<i>n</i>
Wilcoxon Rank Sum	<i>wilcox</i>	<i>m, n</i>
Weibull	<i>weibull</i>	<i>shape, scale = 1</i>

Additional support for generating random numbers is provided by *library(distr)* [11].

A.22 Computing on the Language

The language elements of R are objects, as are data or functions. They can be read and changed like any other data or functions. Chapter 6 of the “R Language Definition” [9] gives details for computing on the language. See also Section 2.1.5, “Function objects” of [9].

<i>Conversions</i>	
<code>parse()</code>	converts input into a list of R expressions. <code>parse</code> executes the parse, but does not evaluate the expression.
<code>deparse()</code>	converts an R expression given in internal representation into a character string.
<code>expression()</code>	generates an R expression in internal representations. <i>Example:</i> <code>integrate <- expression(integral(fun, lims))</code> <i>See also</i> 1.3: mathematical typesetting in plot annotations
<code>substitute()</code>	R expression with evaluation of all defined terms.
<code>bquote()</code>	R expression with selective evaluation. Terms in <code>.()</code> are evaluated. <i>Examples:</i> <code>n<-10; bquote(n^2 == .(n*n))</code>

<i>Evaluation</i>	
<code>eval()</code>	evaluates an expression.

References

- [1] David B. Dahl and contributions from many others. *xtable: Export Tables to LaTeX or HTML*, 2008. R package version 1.5-3.
- [2] D. Firth. Generalized linear models. In D.V. Hinkley, N. Reid, and E.J. Snell, editors, *Statistical Theory and Modeling*, chapter 3, pages 55–82. Chapman and Hall, London, 1991.
- [3] Robert Gentleman and Ross Ihaka. Lexical scope and statistical computing. *Journal of Computational and Graphical Statistics*, 9:491–508, 2000.
- [4] Gabor Grothendieck and Thomas Petzoldt. R help desk: Date and time classes in R. *R News*, 4(1):29–32, June 2004.
- [5] David James and Kurt Hornik. *chron: Chronological Objects Which Can Handle Dates and Times*, 2008. R package version 2.3-24. S original by David James, R port by Kurt Hornik.
- [6] Uwe Ligges and John Fox. R help desk: How can I avoid this loop or make it faster? *R News*, 8(1):46–50, May 2008.
- [7] P. McCullagh and J.A. Nelder. *Generalized linear models*. Number 37 in Monographs on statistics and applied probability. London : Chapman & Hall, 2nd edition, 1989.
- [8] R Development Core Team. *R Data Import/Export*, 2008.
- [9] R Development Core Team. *The R language definition*, 2008.
- [10] R Development Core Team. *Writing R Extensions*, 2008.
- [11] Peter Ruckdeschel, Matthias Kohl, Thomas Stabla, and Florian Camphausen. S4 classes for distributions. *R News*, 6(2):2–6, May 2006.
- [12] William N. Venables and Brian D. Ripley. *S Programming*. Statistics and Computing. Springer, New York, 2000.

