

# Computational Statistics

## Introduction to

---

Günther Sawitzki  
StatLab Heidelberg  
*private version*

June 19, 2008

*in preparation*



---

# Contents

---

<b>Introduction</b>	<b>9</b>
<b>1 Basic Data Analysis</b> <span style="float: right;">1-1</span>	
1.1 R Programming Conventions	1-1
1.2 Generation of Random Numbers and Patterns	1-4
1.2.1 Random Numbers	1-4
1.2.2 Patterns	1-9
1.3 Case Study: Distribution Diagnostics	1-10
1.3.1 Distribution Functions	1-12
1.3.2 Histogram and Related Plots	1-17
Barcharts	1-21
1.3.3 Statistics of Distribution Functions; Kolmogorov-Smirnov Tests	1-22
Monte Carlo Confidence Bands	1-23
1.3.4 Statistics of Histograms and Related Plots; $\chi^2$ -tests	1-30
1.4 Moments and Quantiles	1-35
1.5 R Complements	1-40
1.5.1 Random Numbers	1-40
1.5.2 Graphical Comparisons	1-41
1.5.3 Functions	1-47
1.5.4 Enhancing Graphical Displays	1-51
1.5.5 R Internals	1-54
Parse	1-55
Eval	1-55
Print	1-55

Executing Files	1-56
1.5.6 Packages	1-56
1.6 Statistical Summary	1-58
1.7 Literature and Additional References	1-58
<b>2 Regression</b>	<b>2-1</b>
2.1 General Regression Model	2-1
2.2 Linear Model	2-2
2.2.1 Factors	2-5
2.2.2 Least Squares Estimation	2-6
2.2.3 More Examples for Linear Models	2-16
2.2.4 Model Formulae	2-17
2.2.5 Gauss-Markov Estimator	2-18
2.3 Variance Decomposition and Analysis of Variance	2-20
2.4 Simultaneous Inference	2-27
2.4.1 Scheffé's Confidence Bands	2-27
2.4.2 Tukey's Confidence Intervals	2-29
Case Study: Titre Plates	2-30
2.5 Beyond Linear Regression	2-36
Transformations	2-36
2.5.1 Generalized Linear Models	2-36
2.5.2 Local Regression	2-37
2.6 R Complements	2-41
2.6.1 Discretisation	2-41
2.6.2 External Data	2-41
2.6.3 Testing Software	2-41
2.6.4 R Data Types	2-42
2.6.5 Classes and Polymorphic Functions	2-43
2.6.6 Extractor Functions	2-44
2.7 Statistical Summary	2-45
2.8 Literature and Additional References	2-45

<b>CONTENTS</b>	<b>5</b>
<b>3 Comparison of Distributions</b>	<b>3-1</b>
3.1 Shift/Scale Families, and Stochastic Order	3-4
3.2 <i>QQ</i> plot, <i>PP</i> plot, and Comparison of Distributions	3-6
3.2.1 Kolmogorov Smirnov Tests	3-11
3.3 Tests for Shift Alternatives	3-12
3.4 A Road Map	3-20
3.5 Power and Confidence	3-21
3.5.1 Theoretical Power and Confidence	3-21
3.5.2 Simulated Power and Confidence	3-25
3.5.3 Quantile Estimation by Simulation	3-28
3.6 Qualitative Features of Distributions	3-30
3.7 R Complements	3-31
3.8 Statistical Summary	3-33
3.9 Literature and Additional References	3-33
<b>4 Dimensions 1, 2, 3, ..., <math>\infty</math></b>	<b>4-1</b>
4.1 R Complements	4-2
4.2 Dimensions	4-5
4.3 Selections	4-6
4.4 Projections	4-7
4.4.1 Marginal Distributions and Scatter Plot Matrices	4-8
4.4.2 Projection Pursuit	4-13
4.4.3 Projections for Dimensions 1, 2, 3, ... 7	4-15
4.4.4 Parallel Coordinates	4-16
4.5 Sections, Conditional Distributions and Coplots	4-18
4.6 Transformations and Dimension Reduction	4-24
4.7 Higher Dimensions	4-29
4.7.1 Linear Case	4-29
Partial Residuals and Added Variable Plots	4-29
4.7.2 Non-linear Case	4-31
Example: Cusp Non-linearity	4-31
4.7.3 Melbourne Temperature Data	4-35

4.7.4 “Curse of Dimensionality”	4-38
4.7.5 Case Study: Crabs	4-38
Strategy	4-39
Data Definition	4-39
Data Preparation	4-40
Training Set	4-40
Preliminary Analysis	4-41
Final Analysis	4-42
Discussion	4-42
4.7.6 Case Study: Body Fat	4-43
4.8 High Dimensions	4-55
4.9 Statistical Summary	4-56
4.10 Literature and Additional References	4-56

<b>R as a Programming Language and Environment</b>	<b>A-1</b>
--	------------

A.1 Help and Information	A-1
A.2 Names and Search Paths	A-1
A.3 Customising	A-3
A.4 Basic Data Types	A-4
A.5 Output for Objects	A-6
A.6 Object Inspection	A-7
A.7 System Inspection	A-8
A.8 Complex Data Types	A-9
A.9 Accessing Components	A-11
A.10 Data Manipulation	A-13
A.11 Operators	A-16
A.12 Functions	A-17
A.13 Debugging and Profiling	A-19
A.14 Control Structures	A-21
A.15 Administration and Customisation	A-23
A.16 Input and Output to Data Streams	A-24
A.17 External Data	A-25

<b>CONTENTS</b>	<b>7</b>
A.18 Libraries, Packages	A-27
A.19 Linear Algebra	A-28
A.20 Model Descriptions	A-29
A.21 Graphic Functions	A-31
A.21.1 High Level Graphics	A-31
A.21.2 Low Level Graphics	A-32
A.21.3 Annotations and Legends	A-32
A.21.4 Graphic Parameters and Layout	A-33
A.22 Basic Statistical Functions	A-35
A.23 Distributions, Random Numbers, Densities...	A-36
A.24 Computing on the Language	A-39
<b>Bibliography</b>	<b>Bib-1</b>
<b>Functions and Variables by Topic</b>	<b>Private-1</b>
<b>Function and Variable Index</b>	<b>Private-1</b>
<b>Index</b>	<b>Private-1</b>
<b>B</b>	<b>Private-1</b>
B.1 Spellings and Hyphenations	Private-1
B.2 Formats	Private-1
B.3 Test R function index	Private-1
B.4 R styles	Private-2
B.5 Tables	Private-6
B.6 Help	Private-9
B.7 Indexed Terms	Private-10
<b>C S help — remove for public version</b>	<b>Private-1</b>
C.1 Ch 01	Private-1
<b>D To Do</b>	<b>Private-1</b>
D.1 For release	Private-1
D.2 misc	Private-1



---

# Introduction

---

ch:00

This introduction to R is intended as course material to be used in a compact course or for self-instruction. The course is for students with basic knowledge in stochastics. Notions like distribution function, quantile, expected value and variance are presupposed, as well as some familiarity with statistical features corresponding to these notions. Only a condensed summary is included here. Classical distributions, such as binomial, uniform, and Gaussian should be familiar, together with derived distributions and their asymptotic behaviour. It is not required that the reader has knowledge in statistics properly. But this knowledge is not taught in this course, on the other hand. This course concentrates on the “computing” aspects. Statistical concepts and statistical points of view are introduced and discussed. For an in depth discussion, the reader is referred to statistics courses.

A working knowledge in computer usage is presupposed, at least rudimentary knowledge of programming concepts like variables, loops, functions. No extended knowledge in computing is required.

s:01

## What is R?

R is a programming language, and the name of a software system which implements this language. The R programming language has been developed for statistics and stochastic simulation. By now, it has become a standard in these fields. To be precise, we should use specific terms: the language is called S, its implementation and the environment are called R. The original authors of S are John M. Chambers, R. A. Becker and A. R. Wilks, AT & T Bell Laboratories, Statistics Research Department. The language and its development are documented in a series of books, commonly called by their cover as white ([6]), blue ([2]) and green book ([5]).

For a long time the AT & T implementation of S has been the reference for the S language. Today, S is available in a commercial version called S-Plus <<http://www.insightful.com/>> (based on the AT & T implementation) and as a free software version R<sup>1</sup>, or “Gnu S” <<http://www.r-project.org/>>.



<sup>1</sup> R got its name by accident - the same accident which made the first names of the original authors of R (Ross Ihaka & Robert Gentleman) start with R.

In the meantime, R has become the reference implementation. Essential precisions and - if necessary - even modifications of the language are defined by R. For simplicity, here and in the sequel we use “R language” as a common term even where the precise term should be: the S language using the R implementation.

R is an interpreted programming language. Instructions in R are executed immediately. Besides the original elements of S, R has several extensions. Some of these are introduced in response to recent developments in statistics, some are introduced to open experimental facilities. On the other side, advancements of the S language are taken into account.

The recent (2008) version of R is 2.x. This version is largely compatible with the previous version R 1.x. The essential changes are internal to the system. For initial use, there is no significant difference from R 1.x. For the advanced user, there are three essential innovations:

- *Graphics*: The basic graphics system in R implements a model inspired by pen and paper drawing. A graphics port (paper) is opened, and lines, points or symbols are drawn in this port. In R 2.x there is a second additional graphics system, oriented at a camera/object model. Graphical objects in various positions and orientations are mapped in a visual space.
- *Packages*: The original R had a linear command history and a uniform work space. R 2.x introduced an improved support for “packages” which may have encapsulated information. To support this, language concepts such as “name spaces” and various tools have been introduced.
- *Internationalisation*: Originally, R was based on the English language, and ASCII was the general encoding used. With R 2.x extensive support for other languages and encodings has been introduced. With this, it has become possible to provide localised versions.

Two aspects of R are active areas of recent developments: interactive access, and integration in a networked environment. These and other aspects are part of Omegahat - an attempt to develop a next generation system based on the experiences from R . This more experimental project is accessible at <<http://www.omegahat.org/>>. R does already provide simple possibilities to call functions implemented in other languages like C or FORTRAN. Omegahat extends these possibilities and allows direct access to Java, Perl .... A Java based graphical interface for R is JGR, accessible at <<http://stats.math.uni-augsburg.de/software/>>. A collection of interactive displays for R is in *iplots*, available at the same site.

Recent developments related to R are in <<http://r-forge.r-project.org/>>. Many helpful extensions are in <<http://www.bioconductor.org/>>, a site which is targeted at biocomputing.

s:02 R has been developed for practical work in statistics. Usability often has been given priority over abstract design principles. As a consequence, it is not easy to give a systematic introduction to R. A winding path is chosen here instead: case studies and examples, followed by systematic surveys. Practical work should make use of the rich online material which comes with R. Starting points are the “frequently asked questions” (FAQ)

<<http://www.cran.r-project.org/faqs.html>>. “An Introduction to R” ([19]) is the “official” introduction. This documentation and other manuals can be downloaded from <<http://www.cran.r-project.org/manuals.html>>.

R comes with an extensive online help system providing function descriptions and examples. Once you have become familiar with R, the online help and manuals will become first source of information. For off-line reference, we include some examples of information provided by the help system, with kind permission of The R Foundation for Statistical Computing.

Many R functions are included in the basic system. Other functions must be loaded from libraries. Some of these libraries are distributed with the R system. If additional packages are needed, <<http://www.cran.r-project.org/>> is a first source for downloads, and <<http://r-forge.r-project.org/>> for current projects.

The commercial S-Plus is available in various versions. S-Plus 4.x and S-Plus 2000 use S version 3 and are largely compatible with R. S-Plus 5 is an implementation of S version 4 with some changes which need attention. These programming details of S-Plus are not discussed here. Information about S-Plus is at <<http://www.insightful.com/>>.

### Contents and Outline of this Course

s:03

In its basic version, R contains more than 1500 functions - too many to introduce in just one course, and too many to learn reasonably. The course can only open the way to R.

Course participants can come from various background, with different prerequisites. For pupils and younger students, a mere programming course on technical basics may be appropriate. Later in study, questions about meaningful classification and background will be more important. This is what the present course aims at. The “technical” material gives a skeleton. Beyond this, we try to open the view for statistical questions and to stimulate interest for the background. The course should sharpen the appetite for the substance which may be offered in a subsequent well-founded statistical course.

The first part of this course material is organised by themes, using example topics to illustrate how R can be used to tackle statistical problems.

An appendix gives a collection of R language elements and functions. During the course, it can help as a quick reference and it may serve as a starting point and orientation to access the rich information material which comes bundled with R. After the course, it may serve as a note pad. Finally, on the long run for practical work, the online help and online manuals of R are the prime sources of information.

Using a selection of the exercises, the course can be completed within about four days of work. Conceptually, it is an introduction to statistics with the topic areas

- One sample analysis and distributions
- Regression
- Two sample problems and comparison of distributions

- Multivariate analysis

A generous time slot for exercises is recommended (an additional half day for the introductory exercises, an additional half day for one of the project exercises). With this, the course can be covered in one week, provided follow up facilities are established to answer questions which have come up, and possibilities are available to follow the interest in statistical background which may have resulted.

At a more leisurely pace, chapter 1 with its exercises can be used. This should give a working base to use R, and more material from the subsequent chapters can be added later when it is needed.

Using the course during a term in a weekly class needs more time, since repetitions must be calculated in. Each of the first four chapters will cover about four lectures, plus time for exercises. For this time schedule, a course to cover the statistical background is recommended running in parallel.

For a subsequent self-paced study going into details on R as a programming language, the recommended reading is ([30]).

Statistical literature is evolving, and new publications will be available at the time you read this text. Instead of giving a long list of the relevant literature available at the time this text is written, the sections include keywords which can be used to locate up to date literature.

For economic reasons, part of the illustrations are printed in black and white. Colour versions are available at the web site accompanying this book:

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

Additional material and updates will be available from this site.

### Typographical Conventions

Examples and input code are formatted so that they can be used with “Cut & Paste” and entered as program input. To allow this, punctuation marks are omitted and the input code is shown without a “prompt”. For example

*Example 0.1:*

*Input*

*1 + 2*

*Output*

*3*

may correspond to a screen output of

```
> 1+2
[1] 3
>
```

Depending on the configuration, the prompt ">" may be represented by a different symbol.

## Thanks

Thanks to the R core team for comments and hints. Special thanks to Friedrich Leisch (R core team) and Antony Unwin (Univ. Augsburg). And very special thanks to "ibidem" for comments on the manuscript, for continuing discussion over many years, and for clarifying many views on statistics.

## Literature and Additional References

s:04

**R:Rnotes**  
 [19] R Development Core Team (2000-2007): An Introduction to R.  
 See: <<http://www.r-project.org/manuals.html>>.

**R:Rref**  
 [22] R Development Core Team (2000-2007): R Reference Manual.  
 See: <<http://www.r-project.org/manuals.html>>.

The Omega Development Group (2000): Omega.  
 See: <<http://www.omegahat.org/>>.

**bcw88news**  
 [2] Becker, R.A.; Chambers, J.M.; Wilks, A.R. (1988): The New S Language. NewYork:  
 Chapman and Hall.

**chh92sms**  
 [6] Chambers, J.M.; Hastie, T.J. (eds) (1992): Statistical Models in S. NewYork: Chap-  
 man and Hall.

**clev93vd**  
 [7] Cleveland, W.F. (1993): Visualizing Data. Summit: Hobart Press.

**vr02mass**  
 [31] Venables, W.N.; Ripley, B.D. (2002): Modern Applied Statistics with S. Heidel-  
 berg:Springer.  
 See: <<http://www.stats.ox.ac.uk/pub/MASS4/>>.

**[30]** Venables, W.N.; Ripley, B.D. (2000): Programming in S. Heidelberg:Springer.  
See: <<http://www.stats.ox.ac.uk/pub/MASS3/Sprog>>.

```
$Source: /u/math/j40/cvsroot/lectures/src/SIntro/Rintro/Rnw/S00eintro.Rnw.tex,v $
$Revision: 1.6 $
$Date: 2008/06/18 16:48:41 $
$Name:  $
$Author: j40 $
```

---

## CHAPTER 1

---

# Basic Data Analysis

---

ch:01

### 1.1 R Programming Conventions

s:1.0

Like any programming language, R has certain conventions. Here are the first basic rules.

R Conventions	
Numbers	A point is used as a decimal separator. Numbers can be written in exponential form; the exponential part is introduced by <i>E</i> . Numbers can be complex numbers; the imaginary part is marked by <i>i</i> . <i>Example:</i> 1 2.3 3.4E5 6i+7.8  Numbers can take the values <i>Inf</i> , <i>-Inf</i> , <i>NaN</i> for “not a number” and <i>NA</i> for “not available” = missing. <i>Example:</i> 1/0 results in <i>Inf</i> ; newline 0/0 results in <i>NaN</i> newline <i>NA</i> is used as a placeholder for missing numbers.
Strings	Strings are delimited by " or '. <i>Example:</i> "ABC" 'def' "gh'ij"
Comments	Comments start with # and go to the end of the current line.

To allow for non trivial examples, we anticipate a detail: in R, *a:b* is a sequence of numbers. If  $a \leq b$ , *a:b* is the sequence starting at *a* to at most *b* in steps of 1. If  $a < b$ , *a:b* is the sequence starting at *a* to at least *b* in steps of -1.

<i>R Conventions</i>	
Objects	The basic elements in R are objects. Objects have types, for example <code>logical</code> or <code>integer</code> . Objects can have a class attribute specifying more complex type information. <i>Example:</i> The basic objects in R are vectors.
Names	R objects can have names, by which they can be accessed. Names begin with a letter or a dot, followed by a sequence of letters, digits, or the special characters <code>_</code> or <code>.</code> <i>Examples:</i> <code>x</code> <code>y_1</code> Lower and upper case are treated as different. <i>Examples:</i> <code>Y87</code> <code>y87</code>
Assignments	Assignments have the form <i>Syntax:</i> <code>name &lt;- value</code> or alternatively <code>name = value</code> . <i>Example:</i> <code>a &lt;- 10</code> <code>x &lt;- 1:10</code>
Queries	If only the name of an object is entered, the value of the object is returned. <i>Example:</i> <code>x</code>
Indices	Vector components are accessed by index. The lowest index is 1. <i>Example:</i> <code>x[3]</code> The indices can be specified directly, or using symbolic names or rules. <i>Examples:</i> <code>a[1]</code> the first element <code>x[-3]</code> all elements except the third <code>x[ x^2 &lt; 10 ]</code> all elements where $x^2 < 10$

<i>Help and Inspection</i>	
----------------------------	--

(cont.)→

<b><i>Help and Inspection</i></b> (cont.)	
help	<p>Documentation and additional information about an object can be requested using <code>help</code>.</p> <p>Syntax: <code>help(name)</code></p> <p>Examples: <code>help(exp)</code>  <code>help(x)</code></p> <p>Alternative form <code>?name</code></p> <p>Examples: <code>?exp</code>  <code>?x</code></p> <p>A hypertext (currently HTML) version of R's online documentation is activated by <code>help.start()</code>. This allows search by topics, and a more structured access to information.</p>
Inspection	<p><code>help()</code> can only provide information that has been prepared in advance. <code>str()</code> can inspect the actual state of an object and display this information.</p> <p>Syntax: <code>str(object, ...)</code></p> <p>Examples: <code>str(x)</code></p>

<b>R Conventions</b>	
Functions	<p>Function calls in R have the form</p> <p>Syntax: <code>name(argument ... )</code></p> <p>Example: <code>e_10 &lt;- exp(10)</code></p> <p>This convention even holds when there are no arguments at all.</p> <p>Example: To quit R, you call a “quit” function <code>q()</code>.</p> <p>Function arguments are treated in a very flexible way. They can have default values, which are used if no explicit argument value is given.</p> <p>Examples: <code>log(x, base = exp(1))</code></p>

(cont.)→

<b>R Conventions</b> (cont.)	
	<p>Functions can be <b>polymorphic</b>. For a polymorphic function, the actual function is determined by the class of the actual arguments.</p> <p><i>Examples:</i> <code>plot(x)</code> #a one dimensional serial plot  <code>plot(x, x^2)</code> # a two dimensional scatter plot  <code>summary(x)</code></p>
Operators	<p>When applied to vectors, operators operate on each of the vector components.</p> <p><i>Example:</i> For vectors <code>y, z</code>, the product <code>y*z</code> is the vector of component-wise products.</p> <p>Operators are special functions. They can be called in prefix form (function form).</p> <p><i>Example:</i> <code>"(x, y)"</code></p> <p>When applied to two operands with different lengths, the smaller operand is repeated cyclically.</p> <p><i>Example:</i> <code>(1:2)*(1:6)</code></p>

Our subject are statistical methods. In a first step, we apply the methods in simulations, that is we use synthetical data. Generating these data is largely under our control. This gives us the opportunity to gain experiences with the methods and allows a critical evaluation. Only after that we will use the methods for data analysis.

## 1.2 Generation of Random Numbers and Patterns

s:1.1

### 1.2.1 Random Numbers

Random variables with a uniform distribution can be generated by the function `runif()`. Using `help(runif)` or `?runif` we get information how to use this function:

---

**help(runif)**

---

**Uniform**

*The Uniform Distribution*

---

*Description*

These functions provide information about the uniform distribution on the interval from `min` to `max`. `dunif` gives the density, `punif` gives the distribution function `qunif` gives the quantile function and `runeif` generates random deviates.

*Usage*

```
dunif(x, min=0, max=1, log = FALSE)
punif(q, min=0, max=1, lower.tail = TRUE, log.p = FALSE)
qunif(p, min=0, max=1, lower.tail = TRUE, log.p = FALSE)
runif(n, min=0, max=1)
```

*Arguments*

<code>x, q</code>	vector of quantiles.
<code>p</code>	vector of probabilities.
<code>n</code>	number of observations. If <code>length(n) &gt; 1</code> , the length is taken to be the number required.
<code>min, max</code>	lower and upper limits of the distribution. Must be finite.
<code>log, log.p</code>	logical; if TRUE, probabilities p are given as log(p).
<code>lower.tail</code>	logical; if TRUE (default), probabilities are $P[X \leq x]$ , otherwise, $P[X > x]$ .

*Details*

If `min` or `max` are not specified they assume the default values of 0 and 1 respectively. The uniform distribution has density

$$f(x) = \frac{1}{max - min}$$

for  $min \leq x \leq max$ .

For the case of  $u := min == max$ , the limit case of  $X \equiv u$  is assumed, although there is no density in that case and `dunif` will return `NaN` (the error condition). `runif` will not generate either of the extreme values unless `max = min` or `max-min` is small compared to `min`, and in particular not for the default arguments.

*References*

Becker, R. A., Chambers, J. M. and Wilks, A. R. (1988) *The New S Language*. Wadsworth & Brooks/Cole.

*See Also*

`Random.seed` about random number generation, `rnorm`, etc for other distributions.

*Examples*

```

u <- runif(20)

## The following relations always hold :
punif(u) == u
dunif(u) == 1

var(runif(10000))#- ~ = 1/12 = .08333

```

---

The help information tells us: as an argument for `runif()` we have to supply the number `n` of random variates to generate. As additional arguments for `runif()` we can specify the minimum and the maximum for the range of the random numbers. If we do not specify additional arguments, the default values `min = 0` and `max = 1` are taken. For example, `runif(100)` generates a vector with 100 uniform random numbers with range  $(0, 1)$ . Calling `runif(100, -10, 10)` generates a vector with 100 uniform random numbers in the range  $(-10, 10)$ .

The additional arguments can be supplied in the defined order, or specified by name. If the name of the argument is given, the position can be chosen freely. So instead of `runif(100, -10, 10)` it is possible to use `runif(100, min = -10, max = 10)` or `runif(100, max = 10, min = -10)`. Using the name, it is also possible to set only chosen arguments. For example, if the minimum is not specified, the default value for the minimum is taken: using `runif(100, max = 10)` is equivalent with `runif(100, min = 0, max = 10)`. For better readability, we often write the names of arguments, even if not necessary.

Each execution of `runif()` generates 100 new uniform random numbers. We can store these numbers.

```
x <- runif(100)
```

generates a new vector of random numbers and stores it in the variable `x`.

```
x
```

returns the values. By default, it is written to the output, and we can inspect the result. We get a graphical representation, the *serial plot* - a scatterplot of the entries `x` against its running index, by using

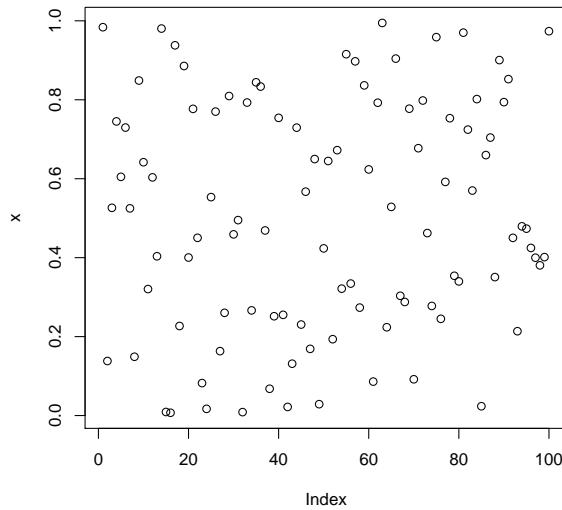
```
plot(x)
```

expl:ch01-Sframe02

***Example 1.1: A Simple Plot***

*Input*

```
x <- runif(100)
plot(x)
```



a:ch01:01-1

Exercise 1.1	
	<p>Try experimenting with these plots and <code>runif()</code>. Do the plots show images of random numbers?</p> <p>To be more precise: do you accept these plots as images of 100 independent realisations of random numbers, distributed uniformly on <math>(0, 1)</math> ?</p> <p>Repeat your experiments and try to note as precisely as possible the arguments you have for or against (uniform) randomness. What is your conclusion?</p> <p>Walk through your arguments and try to draft a test strategy to analyse a sequence of numbers for (uniform) randomness. Try to formulate your strategy as clear as possible.</p> <p style="text-align: right;">(cont.)→</p>

Exercise 1.1	(cont.)
	<p><i>Hint:</i> For comparison, you can keep several plots in a window. The code</p> $\text{par(mfrow} = c(2, 3))$ <p>parametrises the graphics system to show six plots simultaneously, arranged row wise as a <math>2 \times 3</math> matrix (2 rows, 3 columns).</p> <p>The function <code>par</code> is the central function to control graphics parameters. For more information, see <code>help(par)</code>.</p>

**Note:**  
?consider  
simple  
exercises  
for random  
numbers –  
Antony

We reveal the secret<sup>1</sup>: the numbers are not random, but completely deterministic. In the background, `runif()` builds a deterministic sequence  $z_i$  using iterated functions. In the most simple case, for linear congruence generators, subsequent values  $z_i, z_{i+1}$  are generated simply by a linear function. To keep the values in a controlled range, calculation is done modulo an upper bound, that is

$$z_{i+1} = a z_i + b \mod M.$$

The resulting values are re-scaled by

$$\frac{z_i}{M} \cdot (\max - \min) + \min.$$

and returned to us. To start, an initial value  $z_0$ , the `seed`, must be given. Computationally there is just one variable `.Random.seed` which holds the current state  $z_i$  for  $i = 0, 1, \dots$ , and is managed by the system. The successive values returned are in the hand of the user.

The sequence so defined can be regular and soon lead to a periodic solution. With appropriate choice of the parameters however, like in the example given in the footnote, it can result in a long period (in the order of magnitude of  $M$ ) and appear random. But the sequence of numbers is not at all a sequence of independent random numbers, and its distribution is not a uniform distribution on  $(\min, \max)$ .

This is the most simple case. Various other algorithms for random number generation are available, but all follow the same scheme. For more information on available random number generators, see `help(.Random.seed)`. See also Appendix A.23 (page A-36).

Even knowing the secret, a lot of additional knowledge is needed to prove that the generated sequence does not follow the rules which apply to a sequence of independent identically distributed random numbers from a uniform distribution. Sequences which claim to work like random numbers are called **pseudo-random numbers**, if it is important to mark the difference. We use these pseudo-random numbers to generate convenient test data sets. Using these test data sets we can analyse how statistical methods perform under known conditions. In this context, we use pseudo-random numbers as if we had true random numbers.

<sup>1</sup> ... at least part of it. The random number generators used in R can be customised and are usually more complex than the simple variant introduced here. For our discussion, the family of linear congruence generators suffices as an illustration. The usual reference here is the “minimal standard generator” with  $x_{i+1} = (x_i \times 7^5) \bmod 2^{31} - 1$ .

On the other side, we can take pseudo-random numbers as a challenge: are we capable to detect that they are not independent random numbers? If we can detect the difference, we would try to replace the pseudo-random number generator by a better one. But first the challenge goes to us. Are we capable at all to detect that the sequence generated by a linear congruence generator, say, is not random but deterministic? If we cannot: what are the intellectual consequences to draw from this?

### 1.2.2 Patterns

Besides the possibility to generate pseudo-random numbers, R provides several possibilities to generate regular sequences. In many cases these can replace loops, which are common in other languages. Here is a first survey:

R Sequences	
:	generates a sequence from <code>&lt;begin&gt;</code> to at most <code>&lt;end&gt;</code> . <i>Syntax:</i> <code>&lt;begin&gt;:&lt;end&gt;</code> <i>Examples:</i> <code>1:10</code> <code>10.1:1.2</code>
<code>c()</code>	“combine”. combines arguments to a new vector. <i>Syntax:</i> <code>c(..., recursive = FALSE)</code> <i>Examples:</i> <code>c(1, 2, 3)</code> <code>c(x, y)</code> If the arguments are complex data types, the function will descend recursively if called with <code>recursive = TRUE</code> .
<code>seq()</code>	generates general sequences. <i>Syntax:</i> See <code>help(seq)</code>
<code>rep()</code>	repeats argument. <i>Syntax:</i> <code>rep(x, times, ...)</code> <i>Examples:</i> <code>rep(x, 3)</code> <code>rep(1:3, c(2, 3, 1))</code>

a:ch01:01-2

Here “...” denotes a variable list of arguments. We will use this notation frequently.

Exercise 1.2	
	<p>Use</p> $\text{plot}(\sin(1:100))$ <p>to generate a plot of a discretised sine function. Use your strategy from exercise 1.1. Does your strategy detect that the sine function is not a random sequence?</p> <p><i>Hint:</i> If you do not recognise the sine function at first sight, use <code>plot(sin(1:100), type = "l")</code> to connect the points.</p>

**Note:** ?  
show plot  
.RB

Listing the numbers of a data set, as for example the output of a random number generator, rarely helps to detect underlying structures. Simple unspecific graphical representations like the serial plot have some, but only limited use. Even with clear patterns the information provided by these plots is rarely meaningful. Purposeful representations are needed to investigate distribution properties of data. The line plot suggested by the hint in exercise ?? (page ??) already illustrates an analysis beyond scatterplot. It uses a crude interpolation.

### 1.3 Case Study: Distribution Diagnostics

We need specific strategies to detect the presence or the violation of structures. We use random numbers to illustrate how these strategies can look like. Here we concentrate on the distribution properties. Let us assume the sequence consists of independent random numbers with a common distribution. How do we check whether this distribution is a uniform distribution? We will ignore a rescaling to  $(\min, \max)$  - this might be necessary, but is a technical detail which does not affect the investigation substantially. So for now we consider the case  $\min = 0; \max = 1$ .

Realisations of random variables do not allow to read off the distributions directly. This is our critical problem. We need characterisations of the distribution which allow an empirical inspection. Of course we can view the observations as measures. For  $n$  observations  $X_1, \dots, X_n$  we can define the empirical distribution  $P_n$  as  $P_n = \sum(1/n)\delta_{X_i}$ , where  $\delta_{X_i}$  is the point measure at  $X_i$ . Hence

$$P_n(A) = \#\{i : X_i \in A\}/n.$$

**Note:**  
?get damaged/de-  
stroyed  
beyond  
repair

Unfortunately the empirical distribution  $P_n$  of a set of independent observations with common distribution  $P$  is in general very different from  $P$ . Some properties get lost without repair. Infinitesimal properties are among these. So for example  $P_n$  is always concentrated on finitely many points. So empirical distributions are always discrete, irrespective of the underlying distribution. To analyse distributions based on data, we need functionals which can be determined based on realisations of random variable and which can be compared to the corresponding functionals of theoretical distributions. One strategy is to restrict to a family of tests sets that is treatable empirically.

#### S01:ex01 Example 1.1 *Distribution Function*

Instead of the distribution  $P$  we consider its distribution function  $F = F_P$ , where

$$F(x) = P(X \leq x).$$

For an empirical distribution  $P_n$  of  $n$  observations  $X_1, \dots, X_n$ , the corresponding empirical distribution function is

$$F_n(x) = \#\{i : X_i \leq x\}/n.$$

**S01:ex02 Example 1.2 Histogram**

We select disjoint test sets  $A_j, j = 1, \dots, J$ , covering the range of  $X$ . For example for the uniform distribution on  $(0, 1)$  we can chose the intervals

$$A_j = \left( \frac{j-1}{J}, \frac{j}{J} \right]$$

as test sets.

Instead of the distribution  $P$  we consider the vector  $(P(A_j))_{j=1, \dots, J}$ . Its graphical representation is called a **histogram**. The empirical version is the vector  $(P_n(A_j))_{j=1, \dots, J}$ .

We will discuss this example in some detail. Some general lessons can be drawn from this example. We will take several passes, moving from a naive approach to an elaborate statistical approach.

We take the opportunity to point out that the histogram depends critically on the choice of the test sets. In particular, if discretisations in the data meet with an unfortunate choice of the test sets, the results may be misleading. As an alternative to the histogram, we can choose to smooth the data.

**S01:kdens Example 1.3 Smoothing** We replace each data point by a (local) distribution, that is we blur the data points somewhat. We can do this using weight functions. These weight functions are called **kernels** and denoted by  $K$ . We require that the integral over a kernel exists, and conventionally the kernel is normalised so that  $\int K(x) dx = 1$ . Some commonly used kernels are listed in table [T.9](#) and shown in Fig. [T.1](#). For kernels with compact support, the support is chosen to be the interval  $[-1, 1]$  (The R convention is to standardise the kernel, so that the standard deviation is 1.)

By shift and scaling each kernel gives rise to the family

$$\frac{1}{h} K\left(\frac{x - x_0}{h}\right).$$

For kernels, the scale factor  $h$  is called the kernel **bandwidth**. The kernel scaled by  $h$ , is denoted by  $K_h$ :

$$K_h(x) = \frac{1}{h} K\left(\frac{x}{h}\right).$$

The function

$$x \mapsto \frac{1}{n} \sum_i K_h(x - X_i)$$

results in a smoothed display which can replace (or enhance) the histogram.

Kernel	$K(x)$
uniform	$1/2$
triangular	$1 -  x $
Epanechnikov (quadratic)	$3/4(1 - x^2)$
biweight	$15/16(1 - x^2)^2$
triweight	$35/32(1 - x^2)^3$
Gauss	$(2\pi)^{-1/2} \exp(-x^2/2)$

Table 1.9 Some commonly used kernels. See figure 1.1.

ta:ch01kernel

R's density() kernels with bandwidth h = 1

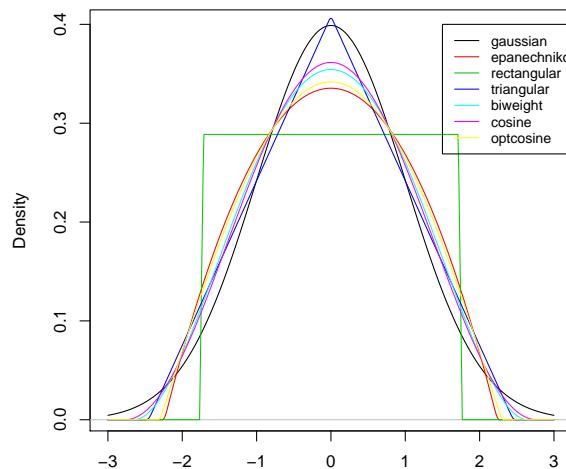


Figure 1.1 Kernels in R. See table 1.9.

fig:Rkernels

For more information, look for the keywords *smoothing* or *kernel density estimation*.

### 1.3.1 First Pass for Example 1.1: Distribution Functions

To test whether a random sequence has a distribution with distribution function  $F$ , we compare  $F$  with the empirical distribution function  $F_n$ . As a special case, we start with a uniform distribution on  $[0, 1]$ . In this special case we have  $F(x) = F_{unif}(x) = x$  for  $0 \leq x \leq 1$ . In general, we have to compute the functions  $F_n$  and  $F$ . A first look tells us:  $F_n$  is a piecewise constant function with jumps at the observations. So we get a complete

subs:s01-vf1

image of  $F_n$  if we evaluate  $F_n$  at the observations  $X_i, i = 1, \dots, n$ . If  $X_{(i)}$  denotes the  $i$ th order statistics, we have  $F_n(X_{(i)}) = i/n$  if  $X_{(1)} < \dots < X_{(n)}$ . We have to compare  $F_n(X_{(i)})$  with the theoretical value  $F(X_{(i)}) = X_{(i)}$ . An R implementation, written with temporary variables, is

```
n <- 100
x <- runif(n)
xsort <- sort(x)
i <- (1:n)
y <- i/n
plot(xsort, y)
```

This is a crude implementation. It does its work in general, but may fail in some cases. If we have ties, that is if we have multiple observations giving the same value, additional precautions have to be taken. The other limitation is that this works for the uniform distribution. In that case, the theoretical distribution function is a diagonal on  $[0, 1]$ , and this is easily compared with the empirical distribution function. For other distributions, we will modify the idea later on. What we still have to solve in general: how do we judge whether a deviation of an empirical distribution function and a theoretical distribution function is significant?

a:ch01:01-tie

Exercise 1.3	Tie Handling
	<p>What is the result if the data set has ties, that is multiple entries with the same value? As an example, look at a data set generated by</p> <pre>x &lt;- runif(n); x &lt;- c(x, x)</pre> <p>How should the correct plot look like?</p>
	<p>Try to modify the implementation so that it behaves correctly in the presence of ties. If you cannot implement it with the programming elements introduced so far: give a list of the programming elements you would need to complete this task.</p>

A more refined implementation would take care of boundary cases or special cases, such as data sets with ties. This sort of details is what in general makes the difference between a first draft and a matured implementation. For now, this implementation serves its purpose.

A line with the theoretical values can be added with

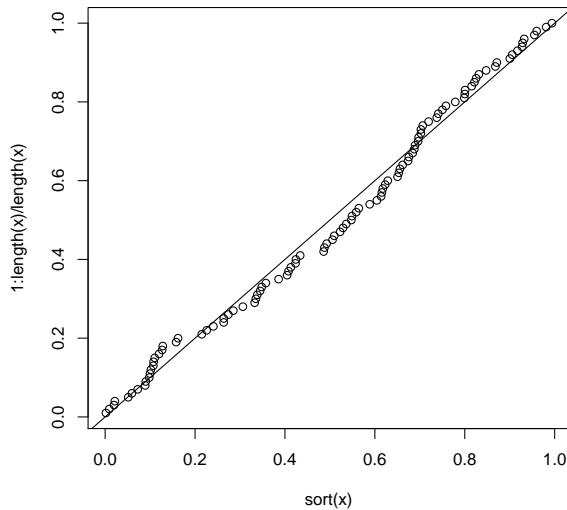
```
abline(0, 1)
```

A more compact implementation using the function `length()` would be:

exp1:ch01-Sframe03

**Example 1.2: Empirical Distribution Function****Input**

```
x <- runif(100)
plot(sort(x), 1:length(x)/length(x))
abline(0, 1)
```



R Functions	
<code>sort()</code>	sorts a vector <i>Example:</i> <code>sort(runif(100))</code>
<code>length()</code>	length of a vector <i>Example:</i> <code>length(x)</code>
<code>abline()</code>	adds a line to a plot <i>Example:</i> <code>abline(a = 0, b = 2)</code>

By default, the `plot()` function adds some annotation. To get a graphic speaking for itself we modify the default annotation to show more specific information. We can replace the default arguments of `plot()`. The argument `main` controls the main title (default: empty). We can add a title, for example

```
plot(sort(x), (1:length(x))/length(x),
     main = "Empirical distribution function\n(X uniform)").
```

Labels for the axis can be controlled using the arguments `xlab` and `ylab`. More information about this and additional arguments can be requested using `help(plot)`. But the help information only provides a forward reference in these cases. The actual information is given by `help(title)`.

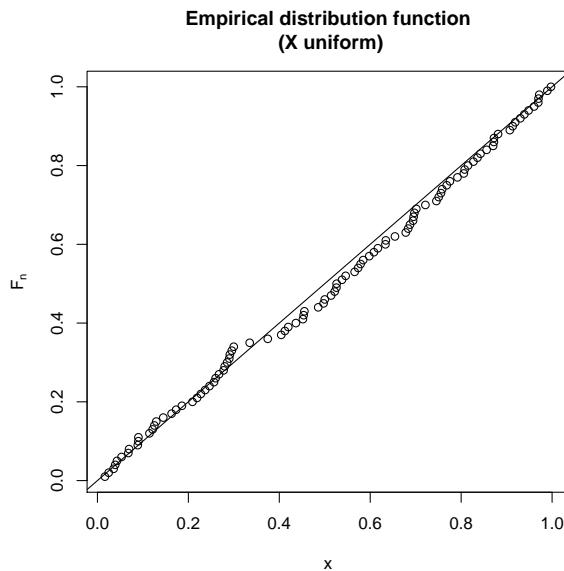
The vertical axis still gives a challenge. Using `ylab = "Fn(x)"` would give a label  $F_n(x)$ . The usual notation however puts the sample size as a subscript, as in  $F_n(x)$ . A hidden property of the levelling functions helps: if a character string is passed as argument, it is displayed “as is”. If an R expression is passed as argument, R tries to display it following the usual mathematical conventions. Details are documented in `help(plotmath)` and examples can be seen using `demo(plotmath)`. To convert a character string to an (unevaluated) R expression, use function `expression()`.

`expl:ch01-1.3`

**Example 1.3: Empirical Distribution Function (Annotated Plot)**

*Input*

```
x <- runif(100)
plot(sort(x), (1:length(x))/length(x),
     xlab = "x", ylab = expression(F[n]),
     main = "Empirical distribution function\n(X uniform)"
)
abline(0, 1)
```



This example serves only for an introduction here. It is not necessary to program distribution functions and their plots from scratch. R already provides a class `ecdf` for empirical distribution functions. If the function `plot()` is applied to an object of class `ecdf`, the

“generic” function `plot()` dispatches internally to the special function `plot.ecdf`, and this draws the distribution function in its specific way. We can reduce our example by just calling `plot(ecdf( runif(100) ))`.

a:ch01:01-ecdf

Exercise 1.4	
	<p>Extend the function call <code>plot(ecdf( runif(10) ))</code> using additional arguments so that the result has the following form:</p>

a:ch01:01-3

Exercise 1.5	Discrimination and Sample Size
	<p>Using <code>norm()</code> you can generate random variables with a Gaussian distribution. Can you tell a Gaussian distribution from a uniform distribution using a serial plot?</p> <p>Use the empirical distribution function</p> <p>Can you tell a Gaussian distribution from a uniform distribution using a plot of the empirical distribution function?</p> <p>Can you tell the sinus series from a uniform distribution using a plot of the empirical distribution function? from a uniform distribution?</p> <p>What is in each case the sample size needed to recognise the differences reliably?</p>

**subs:S01hist1**

### 1.3.2 First Pass for Example S01:ex02 1.2: Histogram

We select test sets  $A_j$ ,  $j = 1, \dots, J$  to cover the range taken by the values of  $X$ . Our strategy: to test whether the data may correspond to a random sequence from a distribution  $P$ , we compare the vector  $(P(A_j))_{j=1,\dots,J}$  with  $(P_n(A_j))_{j=1,\dots,J}$ . For the special case of the uniform distribution on  $(0, 1)$  we can take the intervals

$$A_j = \left( \frac{j-1}{J}, \frac{j}{J} \right]$$

as test sets. Then the vector of theoretical probabilities

$$(P(A_j))_{j=1,\dots,J} = (1/J, \dots, 1/J)$$

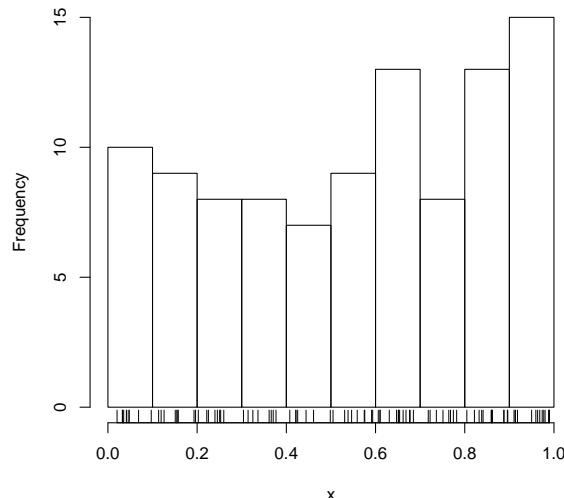
is to be compared to observed relative frequencies  $(\#i : X_i \in A_j)/n$ ,  $j = 1, \dots, J$ . The obvious question in this setting is how to choose the number of cells  $J$ . A preliminary implementation: we use a ready made function to draw histograms. As a side effect, it returns the relative frequencies. The function `rug()` adds the original data to the plot.

**Example 1.4: Histogram With Rug**

*Input*

```
x <- runif(100)
hist(x)
rug(x)
```

Histogram of x



For comparison, we can overlay a kernel density estimation. In contrast to `hist()`, the function `density()` does not draw its result, but prints it. So we have to ask for the

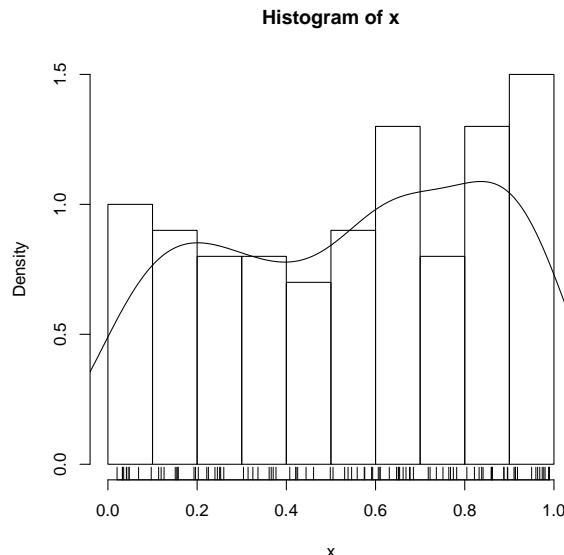
graphical output explicitly. To make scales comparable, we use the argument `probability = TRUE` to get a histogram using a probability scale.

**01-f1histkernel**

**Example 1.5: Augmented Histogram**

*Input*

```
hist(x, probability = TRUE)
rug(x)
lines(density(x))
```



Histograms and kernel density estimators each have their particularities. Histograms suffer from their discretisation, which may interfere with a possible discretisation of the data and give an aliasing effect. Kernel density estimators do not have these discretisation effect, but they blur the data. This may lead to inappropriate boundary effects.

Let us return to the histogram: If we use an assignment

```
whist <- hist(x),
```

the internal information of the histogram is stored as `whist` and be accessed using

```
whist
```

This gives a result like in the following example:

**expl:ch01-f1histkernel**

***Example 1.6: Histogram Data Structure***

<pre>x &lt;- runif(100) whist &lt;- hist(x) whist</pre>	<i>Input</i>
<pre>\$breaks [1] 0.0 0.1 0.2 0.3 0.4 0.5 0.6 0.7 0.8 0.9 1.0  \$counts [1] 6 8 6 17 7 12 10 11 14 9  \$intensities [1] 0.5999999 0.8000000 0.6000000 1.7000000 0.7000000 1.2000000 [7] 1.0000000 1.1000000 1.4000000 0.9000000  \$density [1] 0.5999999 0.8000000 0.6000000 1.7000000 0.7000000 1.2000000 [7] 1.0000000 1.1000000 1.4000000 0.9000000  \$mids [1] 0.05 0.15 0.25 0.35 0.45 0.55 0.65 0.75 0.85 0.95  \$xname [1] "x"  \$equidist [1] TRUE  attr(,"class") [1] "histogram"</pre>	<i>Output</i>

*Counts* gives the cell counts. So these are exactly the numbers we are looking for. The information stored in *whist* consists essentially in five components. Each of these components is a vector. The components have names and can be accessed using these names. So for example

```
whist$counts
```

gives the vector of the cell counts.

<i>Data Structures</i>	
vectors	Components of a vector are accessed by their index. All elements of the vector have the same type. <i>Examples:</i> <code>x</code> <code>x[10]</code>
lists	Lists are complex data structures. The components of a list have names. They can be accessed using these names. Components of a list may be of differing types. <i>Examples:</i> <code>whist</code> <code>whist\$counts</code>

Additional complex data structures are discussed in Appendix [A.8](#) ([S0A1DataStructures](#), [S0C1DataStructures](#))

The choice of histogram cells is data dependent. There are various conventions and recommendations how to choose cells. The details are controlled by arguments for `hist()`. To select test sets of our choice, we have to inspect the calling structure for `hist()`.

`a:ch01:e2hist`

<b>Exercise 1.6</b>	
	<p>Use <code>runif(100)</code> to draw random numbers and generate histograms with 5, 10, 20, 50 cells of equal size. Use repeated samples. Do the histogram plots correspond to what you expect from independent uniform random variates? Try to note your observations in detail.</p> <p>Repeat the experiment with two cells <math>(0, 0.5]</math>, <math>(0.5, 1)</math>.</p> <pre>hist(runif(100), breaks = c(0, 0.5, 1))</pre> <p>Repeat the experiment with random numbers generated by <code>norm(100)</code> and compare the results from <code>runif(100)</code> and <code>norm(100)</code>.</p>

This plot produced by [expl:ch01-fg41histkernel](#) tells only half of the story. As far as the histogram is concerned, the histogram bins can be roughly be guessed from the plot. For the kernel density estimation, the choice of the kernel and the chosen bandwidth are critical, and neither can be read from the resulting display. So the plot has to be enhanced to give this essential information in order to be self-contained.

If you compare `hist` and `density`, you note that R is a developing system. Both functions serve related purposes, but provide different approaches. The obvious difference is that `hist` yields an immediate graphical output while `density` requires an explicit plot request. More subtle differences appear when you try to generate a reproducible analysis. `hist` gives you the details of the histogram used, but does not provide any information on the rules on which this is based. Your have to resume to the history of your commands, and the help information on `hist` (or better the source code of `hist` to be on the sure

side) to retrieve this information. `density` returns information on the call you used as part of the result structure, but the kernel used is not recorded in this information. Again you have to resume to the documentation or the source code.

a:ch01:e3hist

Exercise 1.7	
	<p>Modify example I.5 (page I-18) to include the kernel name and the bandwidth used in the kernel density estimation.            You have to store the result from <code>density()</code> and access its components in analogy to example I.6 (page I-18).</p>

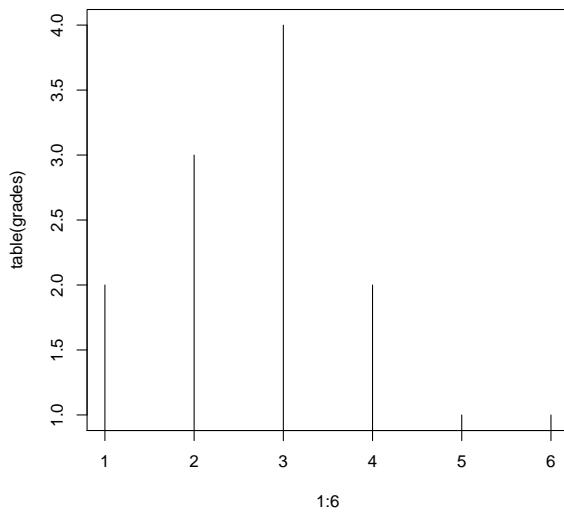
### Barcharts

Note: if the data are not quantitative, but categorical (marked by category labels like “excellent, good, satisfactory, …”, or marked by category number, like “1, 2, 3, …”), a barchart may be more appropriate than a histogram. Simple barcharts are already provided by `plot()`, using the argument `type = 'h'`. To use `plot()`, the raw data must be broken down to frequencies for each level of the variable. This can be done using `table()`.

expl:ch01-bar

#### Example 1.7: Barchart

<pre>grades &lt;- c(2, 1, 3, 4, 2, 2, 3, 5, 1, 3, 4, 3, 6) plot(1:6, table(grades), type = 'h')</pre>	<i>Input</i>
---	--------------



a:ch01:e3bar

Exercise 1.8	Barchart														
	<p>Modify the function call so that the resulting plot is as follows:</p> <table border="1"> <caption>Data for the Grading barchart</caption> <thead> <tr> <th>Grade</th> <th>Count</th> </tr> </thead> <tbody> <tr> <td>1</td> <td>2</td> </tr> <tr> <td>2</td> <td>3</td> </tr> <tr> <td>3</td> <td>4</td> </tr> <tr> <td>4</td> <td>2</td> </tr> <tr> <td>5</td> <td>1</td> </tr> <tr> <td>6</td> <td>1</td> </tr> </tbody> </table>	Grade	Count	1	2	2	3	3	4	4	2	5	1	6	1
Grade	Count														
1	2														
2	3														
3	4														
4	2														
5	1														
6	1														

**ToDo:**

remove  
line caps -  
bar line is  
too long;  
“h” for  
histogram-  
like verti-  
cal lines.

**ToDo:**  
check size  
of barchart  
plot

**ToDo:**  
this is  
thm:so1-beta  
rem. make  
numer-  
tion more  
logical

***1.3.3 Statistics of Distribution Functions; Kolmogorov-Smirnov Tests***

We now move from a naive approach to a statistical point of view. For independent identically distributed variables  $(X_1, \dots, X_n)$  with distribution function  $F$ , in a naive approach we have assumed that  $i/n = F_n(X_{(i)}) \approx F(X_{(i)})$ . We wanted to use this relation to test our distribution assumption. In particular, for a uniform distribution on  $(0, 1)$  this relation reads  $i/n \approx X_{(i)} = F(X_{(i)})$ .

Seen from a statistical point of view, each  $X_{(i)}$  is a random variable. Hence also  $F(X_{(i)})$  is a random variable. Moreover if  $F$  is continuous,  $F(X_{(i)})$  are independent random variables with a uniform distribution on  $[0, 1]$ . We can analyse the distribution of these random variables.

**Theorem 1.4** If  $(X_1, \dots, X_n)$  are independent random variables with common continuous distribution function  $F$ , the  $F(X_{(i)})$  has a beta distribution  $\beta(i, n - i + 1)$ .

*Proof.* → probability theory. Hint: use

$$X_{(i)} \leq x_\alpha \Leftrightarrow (\#j : X_j \leq x_\alpha) \geq i.$$

For continuous distributions, ( $\#j : X_j \leq x_\alpha$ ) has a binomial distribution with parameters  $(n, \alpha)$ .  $\square$

**Corollary 1.5**

$$E(F(X_{(i)})) = i/(n+1).$$

**Exercise 1.9**

Using `help(rbeta)` you get information about the functions available to work with beta distributions. Generate plot for the densities of the beta distribution for  $n = 10, 50, 100$  and  $i = n/4, n/2, 3n/4$ . Use the function `curve()` to generate the plots. For more information, see `help(curve)`.

So on the average, for independent random variables with uniform distribution on  $(0, 1)$  we cannot expect that  $X_{(i)} \approx i/n$ , but on the average we get  $i/(n+1)$ , the expected value of the beta distribution. Hence the reference line should be drawn with `abline(a = 0, b = n/n+1)`.

**Exercise 1.10**

Draw the distribution function with the corrected reference line.

\*

We use the graphical display for a single sample, not for a run of samples. Is the expected value of  $X_{(i)}$  the adequate reference? Are there alternatives which can serve as references?

If you see alternatives, give an implementation.

*Monte Carlo Confidence Bands*

Simulation can also help as to get an impression of the typical fluctuation. We use random numbers to generate a small number of samples, and compare our sample in question with these simulations. For comparison, we generate envelopes of these simulations and check whether our sample lies within the area delimited by the envelopes. If  $x$  is the sample vector to be analysed, with length  $n$ , we use this programming idea:

`expl:ch01-MC06`

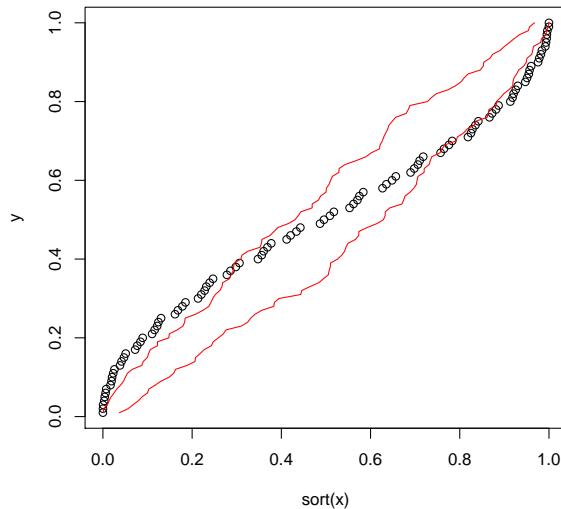
***Example 1.8: Monte Carlo Confidence Bands***

*Input*

```

x <- (sin(1:100)+1)/2
y <- (1:length(x))/length(x)
plot(sort(x), y)
nrsamples <- 19           # nor of simulations
samples <- matrix(runif(length(x)* nrsamples),
nrow = length(x), ncol = nrsamples)
samples <- apply(samples, 2, sort)
envelope <- t(apply(samples, 1, range))
lines(envelope[, 1], y, col = "red")
lines(envelope[, 2], y, col = "red")

```

**ToDo:**

**LB:** meaning?

This example is taken from [\[31\]](#), which is a rich source of ideas for R programming.

A typical R programming strategy is illustrated here. R is an interpreted vector oriented language. Single step interpretation is time extensive. Hence operations with few complex steps can be more efficient than operations using several elementary steps.

- Operations at a vector level are more effective than chains of single elementary operations.
- Iterations and loops should be avoided and replaced by structured vector operations.

a:ch01:e5

Exercise 1.11	
	Make use of <code>help()</code> -function and comment example I.8 step by step. Take a special note of the new functions which are introduced here. <small>exp1:ch01-MC06</small>

**ToDo:**  
ch01: add examples  
f. outer

R Iterators	
<code>apply()</code>	applies a function to the rows or columns of a matrix. <i>Examples:</i> <code>samples &lt;- apply(samples, 2, sort)</code> sorts by column.
<code>outer()</code>	generates a matrix of all pairwise combinations of two vectors, and applies a function to all pairs.

See Appendix A.2 (page A-11) for other iterators provided with R.

For the simulation bands, the hypothesis is that the given sample and the simulations come from the same model. If the curve representing our sample exceeds the limits from our simulation, this indicates that this hypothesis is violated. The approach which is sketched here is called a **Monte Carlo test**. The idea behind this is very general and can be applied in many situations.

a:ch01:e6

Exercise 1.12	
*	<p>Why 19?</p> <p><i>Hint:</i> try to take an abstract simplified view of the problem first: let <math>T</math> be a measurable function and <math>X_0, X_1, \dots, X_{nrsamples}</math> independent samples with a common distribution function.</p> <p>What is <math>P(T(X_0) &gt; T(X_i))</math> for all <math>i &gt; 0</math>.</p> <p>In a second step, give an abstract formulation for the example above. Then specialise for <code>nrsamples = 19</code>.</p>

a:ch01:e7

**ToDo:**  
LB: well, all this depends on the number of data and the number of simulations.

Exercise 1.13	
*	<p>Estimate the coverage probability of the Monte Carlo band by first generating a band as above (how can you draw the band without first making a plot for a special sample?)</p> <p>Next, generate <i>sim</i> simulation samples of uniform random numbers of sample size 100. Count how many simulations give a sample within the band. You have to make your choice of the number <i>sim</i> of simulations (100? 1000? 999?) for this step.</p> <p>Use this information to estimate the coverage probability.</p> <p><i>Hint:</i> <code>any()</code> can be used to evaluate a comparison for a full vector.</p>

**ToDo:**

Monte  
Carlo  
Coverage:  
work out  
as example

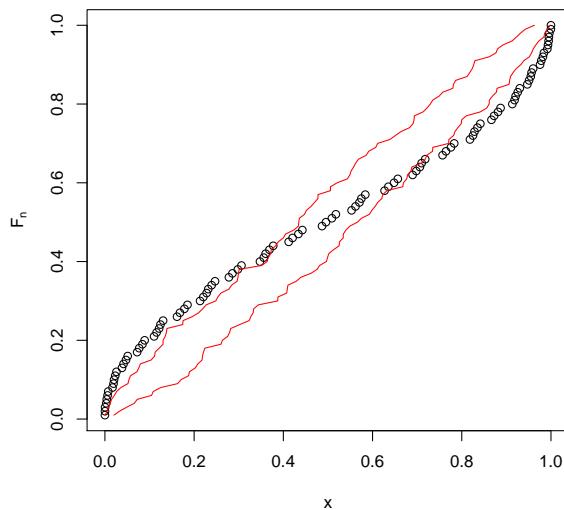
As usual, we want to rework the output so that the resulting plot contains enough information. For the legend, we can proceed as in section II.3.1 (page II-12). The simulation count `nrsamples` needs some thought. If only a fixed count (e.g. 19) is considered, we can use the constant in the legend as usual. If the program fragment should be more versatile, we would need the specific simulation count for each instance. Doing this by hand is error prone, and this possible source of errors can be avoided. The function `bquote()` allows to evaluate a variable, or to calculate it for a specific context. To allow evaluating the current number of simulations, we rearrange the statements so that the number of observations is fixed before `plot()` is called.

expl:ch01-MC06a

**Example 1.9: Monte Carlo Confidence Bands (Augmented)**

<i>Input</i>	<pre> x &lt;- (sin(1:100)+1)/2 y &lt;- (1:length(x))/length(x) nrsamples &lt;- 19 plot(sort(x), y,      main = paste("Monte Carlo Band: ", bquote( .(nrsamples)), " Monte Carlo Samples"),      xlab = 'x', ylab = expression(F[n])) samples &lt;- matrix(data = runif(length(x) * nrsamples),                    nrow = length(x), ncol = nrsamples) samples &lt;- apply(samples, 2, sort) envelope &lt;- t(apply(samples, 1, range)) lines(envelope[, 1], y, col = "red") lines(envelope[, 2], y, col = "red") </pre>
--------------	---

Monte Carlo Band: 19 Monte Carlo Samples



For each simulation new Monte Carlo samples are drawn. So for each invocation we get different Monte Carlo bands, and the bands here are different from those in the previous example.

For practical work it may be necessary to reduce distribution diagnostics to a simple decision problem. This may allow to produce tables or control charts to decide whether a distribution fits to some hypothesis, or to characterise the deviation from a given model. If we want to use tables or simple numbers, additional restrictions are needed: we need to reduce the information contained in the functions  $(F_n, F)$  to simple numbers. One digest is for example the statistic

$$\sup_x |F_n - F|(x).$$

If we want to use this statistic as a criterion, we once again are faced with the task of analysing its distribution.

**Theorem 1.6 (Kolmogorov, Smirnov)** *For a continuous distribution function  $F$ , the distribution of*

$$\sup_x |F_n - F|(x)$$

*is independent of  $F$  (in general, it will depend on  $n$ ).*

*Proof.* → probability theory. For example [9] Lemma 3.3.8. □

**Theorem 1.7 (Kolmogorov)** *For a continuous distribution function  $F$  and  $n \rightarrow \infty$  the statistic*

$$\sqrt{n} \sup |F_n - F|$$

*has asymptotically the distribution function*

$$F_{\text{Kolmogorov-Smirnov}}(y) = \sum_{m \in \mathbb{Z}} (-1)^m e^{-2m^2 y^2} \quad \text{for } y > 0.$$

*Proof.* → probability theory. For example [9] Formula (3.3.11). □

For practical work this means that for a continuous distribution function we can use the following decision strategy: we decide that observations  $(X_1, \dots, X_n)$  do not follow the hypothesis of independent identically distributed observations with distribution function  $F$ , if  $\sup |F_n - F|$  is too large:

$$\sup |F_n - F| > F_{\text{crit}}/\sqrt{n},$$

where  $F_{\text{crit}}$  is taken from the distribution function of the Kolmogorov-Smirnov statistic for sample size  $n$ . This distribution function does not depend on  $F$ . In particular, if we chose the upper  $\alpha$

$F_{\text{crit}} = F_{\text{Kolmogorov-Smirnov}, 1-\alpha}$ , we know that if the hypothesis is valid, this critical value is reached or exceeded at most with probability  $\alpha$ . So we can control the error probability of a false rejection of the hypothesis.

Asymptotically, for large  $n$ , we can use the Kolmogorov-Smirnov approximation instead of the distribution function. If the model distribution  $F$  is not continuous, additional considerations are necessary.

Asymptotics is convenient from a mathematical point of view (because you can use limit theorems), and this is what you find in most text books. But if it comes to practical applications, where  $n = 10$ , or  $n = 100$ , or even  $n = 1000$  or even a larger but always a finite number, asymptotics has always to be taken with caution. What is the finite sample behaviour. In case of the Kolmogorov-Smirnov statistic unfortunately we do not only have asymptotics, but we have bounds which apply to finite sample size and become effective (non-trivial) for very small numbers.<sup>2</sup>

<sup>2</sup> Thanks to L. Birgé for calling attention to this result.

**Theorem 1.8** For all integer  $n$  and any positive  $\lambda$ , we have

$$P(\sqrt{n} \sup |F_n - F| > \lambda) \leq 2e^{-2\lambda^2}.$$

*Proof.* [16] Corollary 2  $\square$

This inequality is valid even if  $F$  is not continuous.

Exercise 1.14	Finite Sample Bounds
	Use the inequality given above to calculate bounds for $\sqrt{n} \sup  F_n - F $ .
	Add finite sample bands to the empirical distribution function.

We want to concentrate on programming issues and do not dive into the details of the Kolmogorov-Smirnov tests. With elementary tools, we can calculate the test statistic  $\sup_x |F_n - F|(x)$  for the uniform distribution. By monotonicity,

$$\sup_x |F_n - F|(x) = \max_{X(i)} |F_n - F|(X(i))$$

and for the uniform distribution

$$\max_{X(i)} |F_n - F|(X(i)) = \max_i |i/n - X(i)|.$$

In R notation, the expression

```
max( abs((1: length(x)) / length(x)) - sort(x)) )
```

is the statistic we are looking for, based on a data vector  $x$ .

Like many commonly used statistics, this statistic and the corresponding distribution function are already implemented in R.<sup>3</sup>

a:ch01:e7ks

Exercise 1.15	
	<p>Using</p> <pre>help(ks.test)</pre> <p>you get information how to invoke the function <code>ks.test</code>.</p> <p>Which results do you expect if you test the following vectors for a uniform distribution:</p> <pre>1:100</pre> <p>(cont.)<math>\rightarrow</math></p>

<sup>3</sup> Different implementations can use other calling structures. The Kolmogorov-Smirnov test is available as `ks.test()`. The distribution function of the test statistic is hidden in the internals of R.

Before R version 2.x however `ks.test()` was not included in the base of R, but was contained in special libraries which had to be loaded explicitly. In R1.x for example the package of classical tests needed to be loaded by `library(ctest)`

Exercise 1.15	(cont.)
	<pre>runif(100) sin(1:100) rnorm(100)?</pre> <p>Perform these tests (scale the values so that they fall into the interval [0, 1], or use a uniform distribution on an interval which is adapted to the data.) and discuss the results.</p>

### 1.3.4 Statistics of Histograms and Related Plots; $\chi^2$ -tests

subs:S01hist2

As we did with the distribution functions, we now move towards a statistical analysis. For simplicity we assume that we have chosen  $A_j, j = 1, \dots, J$  covering the range of  $X$ . The vector of observations  $(X_1, \dots, X_n)$  translates into bin counts  $N_j$

$$N_j = (\#i : X_i \in A_j).$$

If  $(X_i)_{i=1, \dots, n}$  are independent with identical distribution  $P$ ,  $(N_j)_{j=1, \dots, J}$  is a random vector following a multinomial distribution with parameters  $n, (p_j)_{j=1, \dots, J}$  where  $p_j = P(A_j)$ . For the special case  $J = 2$  we have the binomial distribution. Since we are free in our choice of the sets  $A_j$ , we can cover a spectrum of special cases, like

*Median test for symmetry:*

$$A_1 = \{x < x_{0.5}\} \quad A_2 = \{x \geq x_{0.5}\}$$

*Midrange test for concentration:*

$$A_1 = \{x_{0.25} \leq x < x_{0.75}\} \quad A_2 = \{x < x_{0.25} \text{ or } x \geq x_{0.75}\}.$$

For the general case however we must compare the empirical vector of bin counts  $N_j$  with the multinomial distribution, which is not pleasant to calculate. We resort to using approximations. The following approximation goes back to Pearson:

**Lemma 1.9 (Pearson):** For  $(p_j)_{j=1, \dots, J}, p_j > 0$  in the limit  $n \rightarrow \infty$  the following approximation holds:

$$\begin{aligned} P_{mult}(N_1, \dots, N_J; n, p_1, \dots, p_J) \approx \\ (2\pi n)^{-\frac{1}{2}} \left( \prod_{j=1, \dots, J} p_j \right)^{-\frac{1}{2}} \cdot \\ \exp \left( -\frac{1}{2} \sum_{j=1, \dots, J} \frac{(N_j - np_j)^2}{np_j} \right. \\ - \frac{1}{2} \sum_{j=1, \dots, J} \frac{N_j - np_j}{np_j} \\ \left. + \frac{1}{6} \sum_{j=1, \dots, J} \frac{(N_j - np_j)^3}{(np_j)^2} + \dots \right). \end{aligned}$$

*Proof.* → probability theory. For example [14] p. 285.  $\square$

The first term in the exponent of  $e$  is controlled by  $\chi^2 := \sum_{j=1,\dots,J} (N_j - np_j)^2 / np_j$ . This term is called the  $\chi^2$  statistic. At least asymptotically for  $n \rightarrow \infty$  large values of  $\chi^2$  give small probabilities. This motivates to use the  $\chi^2$  statistics approximatively as a measure for goodness of fit. The name comes from the distribution asymptotics:

**Theorem 1.10 (Pearson):** for  $(p_j)_{j=1,\dots,J}, p_j > 0$  in the limit, as  $n \rightarrow \infty$ , the statistics

$$\chi^2 := \sum_{j=1,\dots,J} \frac{(N_j - np_j)^2}{np_j}$$

has a  $\chi^2$  distribution with  $J - 1$  degrees of freedom.

As above, to get a formal decision rule we can fix a critical value  $\chi_{crit}^2$  and reject the hypothesis that the observations  $(X_1, \dots, X_n)$  come from independent identically distributed uniform random variables if the  $\chi^2$  statistic exceeds this value. If we chose the upper  $\alpha$  quantile of the  $\chi^2$  distribution as a critical value, we know that if the hypothesis holds the value  $\chi_{crit}^2$  or a higher value is reached at most with probability  $\alpha$ . So again, at least asymptotically, we can control the error probability of a false rejection of the hypothesis.

The  $\chi^2$  tests are a basic part of R and provided as function `chisq.test()`. As always, the asymptotics used for the “standard”  $\chi^2$  tests has to be taken with caution. Even theoretically, the conditions are more delicate here. It is not sufficient that the sample size is large, but the asymptotics requires that the sample size in each histogram cell must be large. As a second possibility, besides using the asymptotic distribution, `chisq.test()` gives the possibility to calculate compute *p*-values by Monte Carlo simulation (see `help(chisq.test)`).

They are designed to be used for more general “contingency tables”. We just need a special case. In our case, the table is a (one dimensional) vector of counts for the chosen histogram cells. (See also: the R implementation for more general variants in `library(loglin)`.)

a:ch01:e8chisq

Exercise 1.16	
	<p>Use <code>help(chisq.test)</code> to see the calling structure for <math>\chi^2</math> tests. Apply it to test the hypothesis (<math>p_j = 1/J</math>), <math>J = 5</math> on the following vectors of bin counts:</p> $(3 \ 3 \ 3 \ 3 \ 3) \quad (1 \ 2 \ 5 \ 3 \ 3) \quad (0 \ 0 \ 9 \ 0 \ 6).$

a:ch01:e8achisq

Exercise 1.17	
	<p>Which results do you expect if you use a <math>\chi^2</math> test to check the following vectors for a uniform distribution:</p> <pre>1:100 runif(100) sin(1:100) rnorm(100)?</pre> <p>Perform these tests and discuss the results.</p> <p><i>Hint:</i> The function <code>chisq.test()</code> expects a frequency table as input. The function <code>table()</code> can be used to generate a frequency table directly (see <code>help(chisq.test)</code>). But you can also use the function <code>hist()</code> which gives <code>counts</code> as one components as a result.</p>

At first, the approximation for the  $\chi^2$  statistic is only valid, if fixed cells are chosen, independent of the information provided by the sample. Practical histogram algorithms however derive the number of cells and cell bounds based on the sample. Often this implies an estimation of parameters of the distribution. Under some conditions however the  $\chi^2$  asymptotic still holds, as for example illustrated by the following theorem from [24] Section 6b.2:

**To Do:** see  
comment  
LB

**Theorem 1.11** (i) Let the cell probabilities be the specified functions  $\pi_1(\boldsymbol{\theta}), \dots, \pi_k(\boldsymbol{\theta})$  involving  $q$  unknown parameters  $(\theta_1, \dots, \theta_q) = \boldsymbol{\theta}'$ . Further let

- (a)  $\hat{\boldsymbol{\theta}}$  be an efficient estimator of  $\boldsymbol{\theta}$  in the sense of [24] (5c.2.6),
- (b) each  $\pi_i(\boldsymbol{\theta})$  admit continuous partial derivatives of the first order (only) with respect to  $\theta_j$ ,  $j = 1, \dots, q$  or each  $\pi_i(\boldsymbol{\theta})$  be a totally differentiable function of  $\theta_1, \dots, \theta_q$ , and
- (c) the matrix  $M = (\pi_r^{-1/2} \partial \pi_r / \partial \theta_s)$  of order  $(k \times q)$  computed at the true values of  $\boldsymbol{\theta}$  is of rank  $q$ . Then the asymptotic distribution of

$$\chi^2 = \sum \frac{(N_i - n\hat{\pi}_i)^2}{n\hat{\pi}_i} \quad (1.1) \quad \text{eq:chisqrao}$$

is  $\chi^2(k - 1 - q)$ , where  $\hat{\pi}_i = \pi_i(\hat{\boldsymbol{\theta}})$ .

*Proof.* See [24] Section 6b.2.  $\square$

a:ch01:e9

Exercise 1.18	
*	Sketch comparable test environments for fixed and adaptive choice of histogram cells.  (cont.)→

Exercise 1.18	(cont.)
	<p>For fixed and for adaptive choice of histogram cells draw <math>s = 1000</math> samples of size 50 from <code>runif()</code>. Calculate in both settings the formal the <math>\chi^2</math> statistics and plot its distribution functions.</p> <p>Compare the distribution functions.</p>

How to chose the number and breakpoints depends on the purpose of the histogram. One application is to use the envelope curve of the histogram as a density estimator. The built-in algorithms for breakpoint selection of `hist()` are motivated by this application. As a density estimator, the kernel density estimators (Example 1.3) are direct competitors. Kernel density estimators offer more flexibility (for example to control smoothness) whereas in general histograms are more easy to calculate. More often, histograms are used to detect **data features**. However they are unless they are tailored for special cases, more efficient plots are available for this purpose. We will return to the topic of data features in section 3.6.

### ***Repeated Samples***

So far, we have concentrated on inspecting the distribution of random variables. We can continue this inspection. If  $(X_1, \dots, X_n)$  are independent random numbers with identical uniform distribution, for fixed cells the  $\chi^2$  statistics is approximatively distributed as  $\chi^2$ , and  $\kappa := \sqrt{n} \sup |F_n - F|$  has the Kolmogorov-Smirnov distribution.

We can take repeated samples  $(X_{1,j}, \dots, X_{n,j})_{j=1..m}$  and from these we can calculate statistics  $\chi^2_j$  and  $\kappa_j$ . Given independent, identically distributed random variables to start with, the statistics have a  $\chi^2$  resp. a Kolmogorov-Smirnov distribution. Taking these repeated samples we can not only inspect the distribution of the single observations, but the joint distribution of samples, each consisting of  $n$  sample elements.

Exercise 1.19	
	<p>For <math>n = 10, 50, 100</math>, draw 300 samples using <code>runif(n)</code>. For each sample, calculate the <math>\chi^2</math> and the Kolmogorov-Smirnov statistic.</p> <p>You have to chose a <math>\chi^2</math> test to take. What is your choice?</p> <p>Plot the distribution functions of these statistics and compare them to the theoretical (asymptotic) distributions.</p> <p>Are there any indications against the assumption of independent uniform random numbers?</p> <p><i>Hint:</i> the functions for the <math>\chi^2</math> and Kolmogorov-Smirnov test keep their internal information as a list. To get the names of the list elements, you can create sample object. For example, use  <code>names(chisq.test(runif(100)))</code>.</p>

***Power***

So far in our discussion the uniform distribution the target model distribution, our “hypothesis”. We have discussed how various methods will behave, if this hypothesis holds. The distribution properties derived from the model assumption can indicate how to find critical values for formal tests. We reject the hypothesis, if the observed test statistics are too extreme. What “too extreme” means is determined by the distributions we derived from the hypothesis. This leads to decision rules like:

reject the hypothesis, if  $F_{\chi^2}(\chi^2) \geq 1 - \alpha$

or

reject the hypothesis, if  $F_{Kolmogorov-Smirnov}(\kappa) \geq 1 - \alpha$

for chosen (small) values of  $\alpha$ .

In the next step, if we have fixed a formal decision rule, we can ask for the power of this rule to come to a rejection if the hypothesis actually does not hold. A more precise analysis of the power of decision rules is a theme in classical lectures in statistics. The possibilities discussed so far allow us to analyse the performance using Monte Carlo strategies.

As framework for simulations we chose a family of alternatives. The uniform distribution blends into the family of beta distributions with densities

$$p_{a,b}(x) = \frac{\Gamma(a+b)}{\Gamma(a)\Gamma(b)} x^{a-1} (1-x)^{b-1} \quad \text{for } a > 0, b > 0 \text{ and } 0 < x < 1.$$

We chose alternative distributions from this family. From these we draw repeated samples, and apply our decision rule formally on each of these.

We note whether the procedure leads to a rejection of the hypothesis or not. Given a choice of the sample size  $n$  and a number of simulations  $m$ , and given a choice of a limiting probability  $\alpha$  we get a table

$$(a, b) \mapsto \# \text{ simulations leading to a rejection of the hypothesis.}$$

In particular for the uniform distribution  $(a, b) = (1, 1)$  we expect approximately  $m \cdot \alpha$  rejections. For distributions other than the uniform hypothesis, a decision procedure is more powerful if the proportion of rejections is higher.

a:ch01:e11ts

Exercise 1.20	
**	<p>Analyse the power the Kolmogorov-Smirnov test and the <math>\chi^2</math> tests. Select values for <math>n, m</math> and <math>\alpha</math>, and chose 9 pairs for <math>(a, b)</math>. What are your arguments for your choices?</p> <p>Use your chosen parameters to draw samples from <code>rbeta()</code>.</p> <p>Apply the Kolmogorov-Smirnov test and a <math>\chi^2</math> test with 10 cells of equal size on <math>(0, 1)</math>.</p>
	(cont.)→

Exercise 1.20	(cont.)
	<p>Chose alternative parameters <math>(a, b)</math> so that you can compare the decision rules along the following lines:</p> <ul style="list-style-type: none"> <li>i) <math>a = b</math></li> <li>ii) <math>b = 1</math></li> <li>iii) <math>a = 1</math></li> </ul> <p>and run these simulations.</p>
	<p>Chose alternative parameters <math>(a, b)</math> so that you can compare the decision rules over the range <math>0 &lt; a, b &lt; 5</math>.</p> <p>Your conclusions?</p> <p><i>Hint: outer(x, y, fun)</i> applies a function <i>fun()</i> to all pairs of values from <i>x, y</i> and returns the result as a matrix.</p> <p>Using</p> $\text{contour}()$ <p>you can generate a contour plot. See <i>demo("graphic")</i>.</p>

a:ch01:e12pseudo

Exercise 1.21	
**	<p>Design a test strategy to unmask “pseudo-random numbers” .</p> <p>Test this strategy using simple examples</p> <ul style="list-style-type: none"> <li>i) <math>x \quad x = 1..100 \mod m</math> for convenient <math>m</math></li> <li>ii) <math>\sin(x) \quad x = 1..100</math></li> <li>iii) ...</li> </ul> <p>Do you tag these sequences as “not random”?</p> <p>Now try to unmask the random number generators provided by R. Can you identify the generated sequences as “not random”?</p>

## 1.4 Moments and Quantiles

Distribution functions or densities are no easy objects from a mathematical point of view. In general, the space of functions is infinite dimensional. Finite dimensional geometrical or optimisation arguments cannot be applied directly. To simplify analysis one often resorts to finite dimensional descriptions.

Historically, moments have plaid an important role: probability distributions are viewed as mass distributions, and the moments are defined in analogy to the moments in mechanics. The first moment, corresponding to the center of gravity, is called ***expected value*** in statistics.

**Definition 1.12** If  $X$  is a real valued random variable with distribution  $P$ , the expected value of  $X$  is defined as

$$E_P(X) := E(X) := \int X dP$$

provided this integral exists.

The second moment and higher moments are usually centred. For the second (central) moment, the **variance**, we have the following definition:

**Definition 1.13** If  $X$  is a real valued random variable with distribution  $P$ , its variance is defined as

$$\text{Var}_P(X) := \text{Var}(X) := \int (X - E(X))^2 dP.$$

provided this integral exists.

The integral does not always give defined values. Hence the moments do not always exist. But if they exist, they give a first information about the distribution. The expected value is often interpreted as the “statistical average” or mean of the distribution. The root of the variance, the **standard deviation**, is interpreted as the “statistical spread”.

The definitions can be applied to empirical distributions as well. This gives a first possibility to estimate the moments of an unknown theoretical distribution from the data. For the expected value, the mean, we have an unbiased estimator:

$$E_P(E_{P_n}(X)) = E_P(X),$$

That is in the statistical average the empirical mean and the mean of the underlying theoretical distribution coincide (provided they are defined).

For the variance the unbiasedness does not hold, but we have

$$\frac{n}{n-1} E_P(\text{Var}_{P_n}(X)) = \text{Var}_P(X),$$

for  $n > 1$ . The mathematical background is that the expected value is a linear operator. It commutes with linear operators. But the variance is a quadratic operator, and that requires a correction if we want an unbiased estimator. The variance with this bias correction is often called the **empirical variance** or **sample variance**.

For the estimation of the first two moments of a vector of random numbers, R provides two functions: `mean()` estimates the mean and `var()` estimates the variance, using the sample variance. The function `sd()` gives the standard deviation of a vector.

A quite different approach is to reduce the distribution function to a small set of percentage points, or **quantiles**, such as the median (the 50% point) and the lower and upper quartile (the 25% and 75% point). With this approach, for example the **interquartile range** (the difference between lower and upper quartile) can be used as a measure of spread. Since distribution functions can have steps and flat parts, an exact definition of quantiles has to be more elaborate:

**Definition 1.14**  $x \in \mathbb{R}$  is a  $p$  **quantile** of  $X$ , if  $F_X(x) \geq p$  and  $F_X(x') \leq p$  for all  $x' < x$ .

Exercise 1.22	
	<p>Generate a sample of random variables with sample size 100 from the distributions with the following densities:</p> $p(x) = \begin{cases} 0 & x < 0 \\ 1 & 0 \leq x \leq 1 \\ 0 & x > 1 \end{cases}$ <p>and</p> $p(x) = \begin{cases} 0 & x \leq 0 \\ 2 & 0 < x \leq 1/4 \\ 0 & 1/4 < x \leq 3/4 \\ 2 & 3/4 < x \leq 1 \\ 0 & x > 1 \end{cases}$ <p>Estimate the mean, variance and standard deviation in each of these.</p> <p>Repeat the estimation for 1000 samples. Analyse the distribution of estimated mean, variance and standard deviation for repeated samples.</p>

Moments can be calculated using simple arithmetic operations. Their combination (exact or approximate) follows simple laws. However they are quite sensitive. Even shifting minimal probability mass can lead to a breakdown. For the empirical distribution this means: if a proportion of  $1 - \varepsilon$  of the observed data follows a model distribution and a proportion of  $\varepsilon$  comes from some different distribution, the moments can take any value, even for arbitrarily small values of  $\varepsilon$ . Quantiles are much more robust against a breakdown. At least 50% of the data have to be “outliers” to lead to a breakdown of the median, while changes in just one data point can give any value to the mean.

With the availability of sufficient computing power, quantiles have gained more importance as descriptors. Calculating quantiles implicitly requires sorting. Hence it is more complex than the calculation of moments. Rules for combination of quantiles are not as simple as those for moments and often need explicit calculations. But with the technical resources available now this is not an essential restriction.

R provides several functions to work with quantiles.

`quantile()` is an elementary function to calculate quantiles. The function `summary()` gives a summary of the distribution information, which includes information about quantiles.

a:ch01:e19

Exercise 1.23	
	<p>Generate a sample of 100 random variables from the distributions of exercise 1.22.</p> <p>Estimate the median, and the lower and upper quartile.</p> <p>Repeat the estimation for 1000 samples. Analyse the distribution of the estimated median, lower and upper quartile from repeated samples.</p>

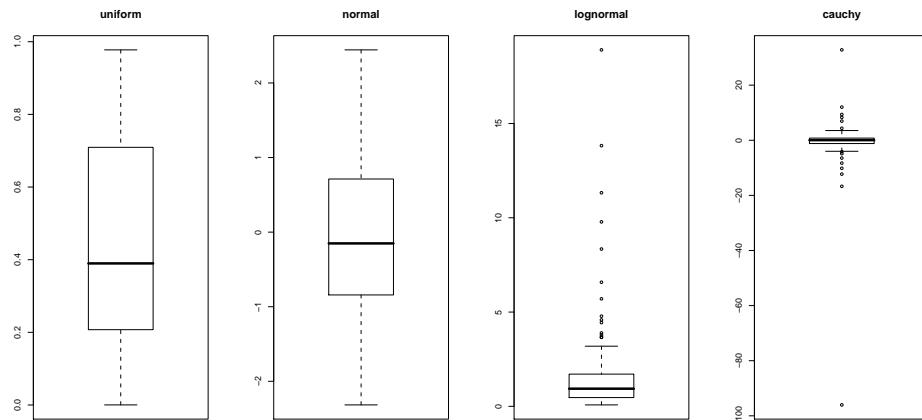
`boxplot()` gives a graphical representation of a quantile based summary. The “box&whisker plot” which is used here has many variations. So for interpretation of a box&whisker plot it is necessary to get detailed information which convention is used. Usually a box is used to mark the central part of the distribution. In the standard version a line is used to mark the median, and a “Box” goes from the median of the upper half to the median of the lower half. Roughly this corresponds to the upper and lower quartile. The finer definition takes care that the information is still reliable if the data contains ties, that is repeated observations with the same values. The “whisker” describes the adjacent area. Outliers are marked in particular.

expl:ch01-box

#### Example 1.10: Box-and-Whisker Plot

Input

```
oldpar <- par(mfrow = c(1, 4))
boxplot(runif(100), main = "uniform")
boxplot(rnorm(100), main = "normal")
boxplot(exp(rnorm(100)), main = "lognormal")
boxplot(rcauchy(100), main = "cauchy")
par(oldpar)
```



a:ch01-box

**To Do:**  
comment  
differences

Exercise 1.24	
	Modify example <a href="#">expl:ch01-box</a> so that the plots are comparable: adjust the location so that the medians are on the same height. Adjust the scales so that the inter quartile ranges have same length.

Theorem [I.4](#) gives a possibility to determine confidence intervals for quantiles that are valid, independent of the the underlying distribution.

To get an upper bound for the  $p$ -quantile  $x_p$  of a continuous distribution function by an order statistics  $X_{(k:n)}$  with confidence level  $1 - \alpha$  we look for

$$\min_k : P(X_{(k:n)} \geq x_p) \geq 1 - \alpha.$$

But  $X_{(k:n)} \geq x_p \iff F(X_{(k:n)}) \geq p$  and by theorem [I.4](#) we have

$$P(X_{(k:n)} \geq x_p) = 1 - F_{beta}(p; k, n - k + 1).$$

So we can determine  $\min_k$  directly using the beta distribution, or we use the relation with the binomial distribution and calculate  $k$  as

$$\min_k : P_{bin}(X \leq k - 1; n, p) \geq 1 - \alpha.$$

a:ch01:e4

Exercise 1.25	
	For continuous distributions and the median $X_{med}$ we have $P(X_i \geq X_{med}) = 0.5$ . Hence we can find a $k$ such that $k = \min\{k : P(X_{(k)} \leq X_{med}) < \alpha\}$ and $X_{(k)}$ as an upper bound for the median with confidence level $1 - \alpha$ . Use this idea to construct a confidence interval for the median with confidence level $1 - \alpha = 0.9$ .
	Modify the box & whisker plot to show this interval. <i>Hint:</i> You need the distribution function $F_X$ , evaluated at the position marked by the order statistics $X_{(k)}$ . The distributions of $F_X(X_{(k)})$ is discussed in <a href="#">I.4</a> .
	The box & whisker plot offers an option <code>notch = TRUE</code> to mark confidence intervals. Try to use the documentation to find out how a <code>notch</code> is calculated. Compare your confidence intervals with those marked using <code>notch</code> .
	(cont.)→

Exercise 1.25		(cont.)
*		Use an analogous strategy to get a distribution independent confidence interval for the inter quartile distance.
***		Augment the box & whisker plot so that it gives information about the scale in a way which is statistically reliable. <i>Hint:</i> Why is it not sufficient to mark confidence intervals for the quartiles?

## 1.5 R Complements

s:1.3

### 1.5.1 Random Numbers

If we had independent identical uniform distributed random numbers, we could generate random numbers with many other distributions. For example:

**Lemma 1.15** (*Inversion method: If  $(U_i)$  is a sequence of independent identically  $U[0, 1]$  distributed random variables and  $F$  a distribution function, then  $(X_i) := (F^{-1}U_i)$  is a sequence of independent identically distributed random variables with distribution  $F$ .*)

From an analytical point of view, this lemma is only usable if  $F^{-1}$  is known. Numerically, it has a much wider use. Instead of  $F^{-1}$  we often use only approximations, sometimes just an inversion table.

The inversion method is a method to use uniformly distributed random numbers to get random numbers with other distributions of interest. Other, some times much more effective methods to get random numbers with specific distributions are discussed in the literature on statistical simulation.

#### ToDo:

www: add references

For a range of distributions, transformed random number generators are supplied with R. A list is given in appendix (page [A-36](#)). For each distribution family there is a series of functions. The names for these functions is derived from the short name of the distribution. For the family  $xyz$ ,  $rxyz$  is a function that generates random numbers.  $dxyz$  calculates the density or the point measure for this family.  $pxyz$  gives the distribution function and  $qxyz$  the quantiles<sup>4</sup>.

Some [selected distributions](#). For more distributions distributions see [A.23](#) (page [A-36](#)).

Distribution	Random Numbers	Density	Distribution Function	Quantiles
Binomial	<code>rbinom</code>	<code>dbinom</code>	<code>pbinom</code>	<code>qbinom</code>

(cont.)→

<sup>4</sup> Confusingly, with the notations which are common in statistics, we have  $p_{xyz} \equiv dxyz$  and  $F_{xyz} \equiv pxyz$ .

<i>Distribution</i>	<i>Random Numbers</i>	<i>Density</i>	<i>Distribution Function</i>	<i>Quantiles</i>
Hypergeometric	<code>rhyper</code>	<code>dhyper</code>	<code>phyper</code>	<code>qhyper</code>
Poisson	<code>rpois</code>	<code>dpois</code>	<code>ppois</code>	<code>qpois</code>
Gauss (normal)	<code>rnorm</code>	<code>dnorm</code>	<code>pnorm</code>	<code>qnorm</code>
Exponential	<code>rexp</code>	<code>dexp</code>	<code>pexp</code>	<code>qexp</code>

### 1.5.2 Graphical Comparisons

Differences between simple geometrical forms are perceived easier than differences between general graphs of similar form. So it can be helpful to choose representations leading to simple forms such as straight lines. For example, to compare two distribution functions  $F, G$ , instead of comparing the function graphs we can consider the graph of

$$x \mapsto (F(x), G(x)).$$

This graph is called the *PP plot* or *probability plot*. If the distributions coincide, this graph shows a diagonal straight line. Deviations from the diagonal are perceived easily.

Alternatively, we can use the observation scale as a reference. So we consider the graph of

$$p \mapsto (F^{-1}(p), G^{-1}(p))$$

This graph is called the *QQ plot* or *quantile plot*. If the distributions coincide, this graph again shows a diagonal straight line.

For the special case of uniform distributions on  $[0, 1]$  on this interval we have  $x = F(x) = F^{-1}(x)$ , that is the *QQ* plot and the *PP* plot coincide. In this case, they both are just the graph of the distribution function. For non uniform distributions in the *PP* plot the distribution functions are standardised to the probability scale  $[0, 1]$ , and in the *QQ* plot they are re-scaled to the observation scale.

**ToDo:**

more precise word  
than easier

a:ch01qqpp

Exercise 1.26	
	Generate a <i>PP</i> plot of the $t(\nu)$ distribution against the standard normal distribution in the range $0.01 \leq p \leq 0.99$ for $\nu = 1, 2, 3, \dots$
	Generate a <i>QQ</i> plot the $t(\nu)$ distribution against the standard normal distribution in the range $-3 \leq x \leq 3$ for $\nu = 1, 2, 3, \dots$
	How large must $\nu$ be so that the $t$ distribution is barely different from the normal distribution in these plots?
	How large must $\nu$ be so that the $t$ distribution is barely different from the normal distribution if you compare the graphs of the distribution functions?

If the distributions are related by an affine transformation in the observation space, the  $QQ$  plot still shows a straight line; slope and intercept represent the affine transformation. This applies for example to the family of normal distributions: if  $F$  is the standard normal distribution  $N(0, 1)$  and  $G = N(\mu, \sigma^2)$ , the  $QQ$  plot is a straight line with intercept  $\mu$  and slope  $\sigma$ .

For the empirical distributions, corollary I.5 applies: instead of  $i/n$  we use a reference point which corrects for the skewness of the distribution, so that on the average a straight line is generated. The quantile plot, using this correction for empirical distributions, is provided as function `qqplot()`. For the special case of the normal distribution there is a variation of `qqplot()` available as `qqnorm()` to compare an empirical distribution with the theoretical normal distribution.

By transforming to probability scale or observation scale the graphical procedures gain power. So for example using a sample size of  $n = 50$ , to tell the distribution function of the normal distribution from that of a uniform distribution often needs a well trained observer. In the normal  $QQ$  plot however the uniform samples show up as markedly non-linear, whereas normal sample largely show linear pictures, and the difference becomes obvious.

To illustrate this, we generate random samples:

---

*Input*

```
unif50 <- runif(50)
unif100 <- runif(100)
norm50 <- rnorm(50)
norm100 <- rnorm(100)
lognorm50 <- exp(rnorm(50))
lognorm100 <- exp( rnorm(100))
```

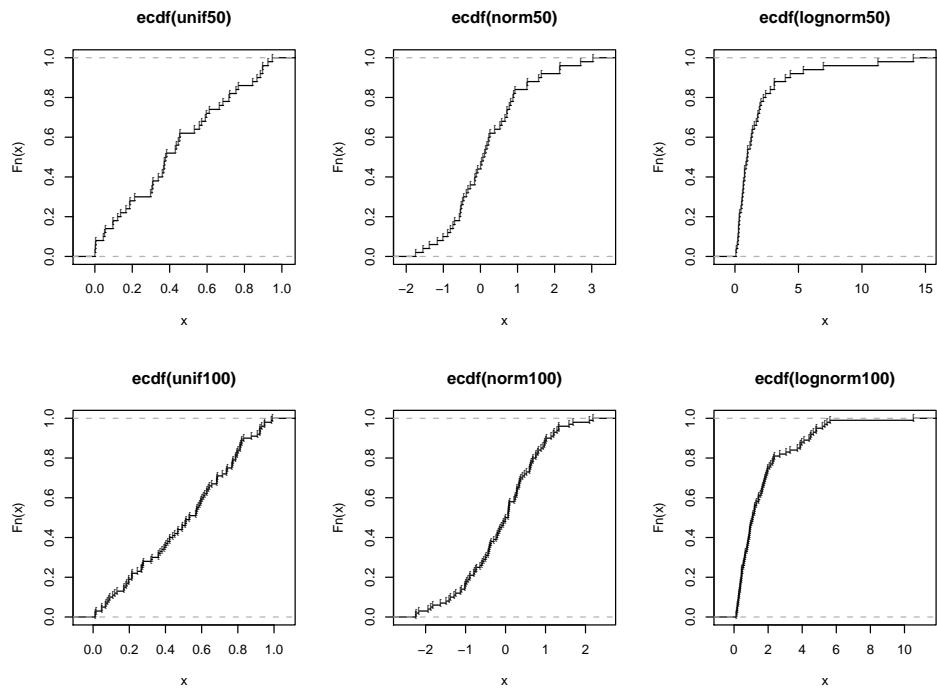
Using these data sets, we generate plots of the distribution functions.

**ToDo:**

Add comment:  
these are tools, not recipes

**Example 1.11: Distribution Functions (Various Distributions)**Input

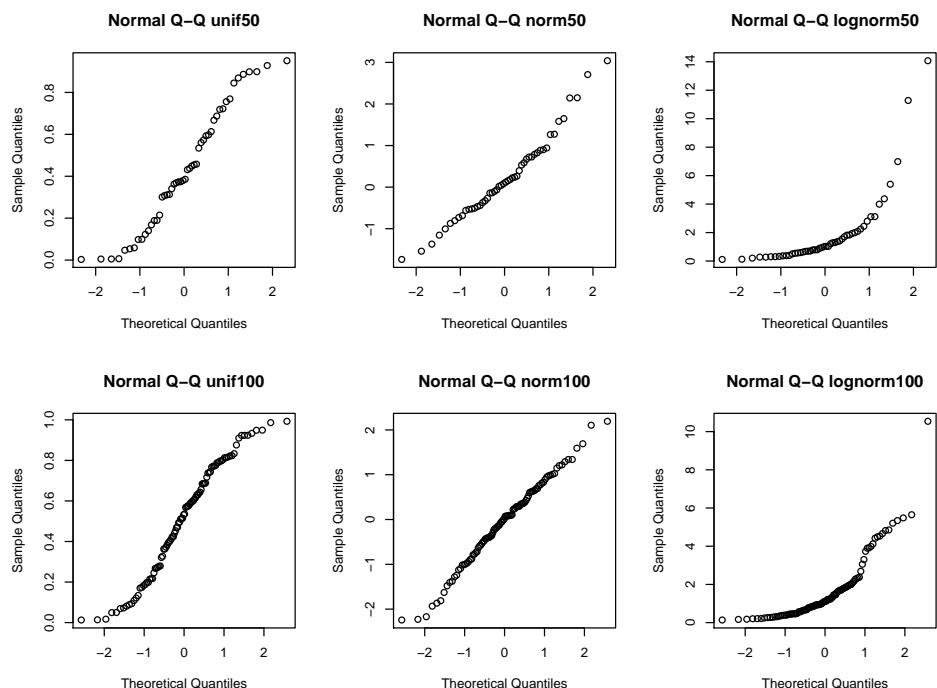
```
oldpar <- par(mfrow = c(2, 3))
plot(ecdf(unif50), pch = "[")
plot(ecdf(norm50), pch = "[")
plot(ecdf(lognorm50), pch = "[")
plot(ecdf(unif100), pch = "[")
plot(ecdf(norm100), pch = "[")
plot(ecdf(lognorm100), pch = "[")
par(oldpar)
```



expl:ch01-compqq For comparison, here are the corresponding normal QQ plots for the same data:

**Example 1.12: Normal QQ Plots (Various Distributions)**Input

```
oldpar <- par(mfrow = c(2, 3))
qqnorm(unif50, main = "Normal Q-Q unif50")
qqnorm(norm50, main = "Normal Q-Q norm50")
qqnorm(lognorm50, main = "Normal Q-Q lognorm50")
qqnorm(unif100, main = "Normal Q-Q unif100")
qqnorm(norm100, main = "Normal Q-Q norm100")
qqnorm(lognorm100, main = "Normal Q-Q lognorm100")
par(oldpar)
```



a:ch01:e1vf

<b>Exercise 1.27</b>	
	Use <i>PP</i> plots instead of distribution functions to illustrate the $\chi^2$ - and Kolmogorov-Smirnov approximations.

a:~~ch01~~  
e2qq  
better  
RExercise

Exercise 1.28	
	Use <i>QQ</i> plots instead of distribution functions. Can you add confidence regions to these plots with the help of the $\chi^2$ - resp. Kolmogorov-Smirnov statistics?

To get an impression of the fluctuation, we have to compare the empirical plots with typical plots of a model distribution. A plot matrix is a simple way. Here is an example for the normal *QQ* plot, implemented as a function:

---

```
Input
qqnormx <- function(x, nrow = 5, ncol = 5, main = deparse(substitute(x))){
  oldpar <- par(mfrow = c(nrow, ncol))
  qqnorm(x, main = main)
  for (i in 1:(nrow*ncol-1)) qqnorm(rnorm(length(x)), main = "N(0, 1)")
  par(oldpar)
}
```

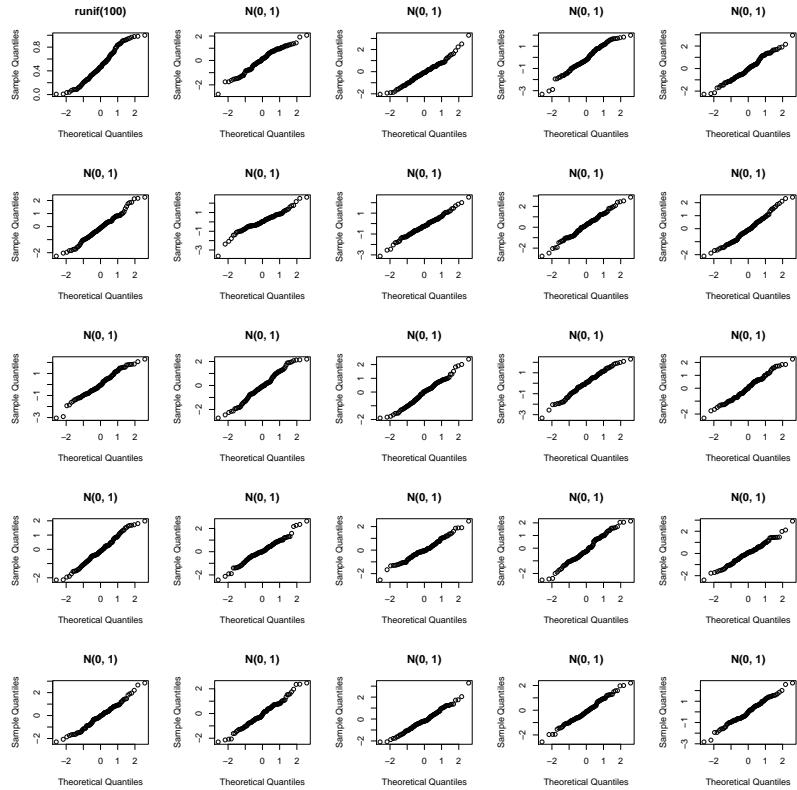
---

In this example we used a *for* loop. Like all programming languages, R has control structures such as loops or conditional statements. A summary of control structures in R is in appendix [A.14](#).

Exercise 1.29	
	Generate a matrix of dimensions $(nrow * ncol - 1), length(x)$ with random numbers and use <i>apply()</i> to avoid the loop. <i>Hint:</i> see example <a href="#">I.8</a> (page I-23).

Deviations from a linear structure should be considered as fluctuation if they stay within the frame of the simulated examples. If the data set under investigation is too extreme in comparison with the simulated examples, this indicates a contradiction to the model assumptions.

[expl:ch01-qqmat](#)

**Example 1.13: Graphical Monte Carlo Test****Input**`qqnormx(runif(100))`

On the long run it is worth to modify the plot functions so that they give information about the fluctuation which is to be expected. In example 1.9 we have constructed Monte Carlo bands for the distribution function. We can generalise this idea for the *PP* plot and the *QQ* plot. All that is needed is to transform the bands to the adequate scale used in the plot.

`a:ch01:e3ppqq`

<b>Exercise 1.30</b>	
	Use <code>rnorm()</code> to generate pseudo-random numbers for the normal distribution for sample size $n = 10, 20, 50, 100$ . For each sample, generate a <i>PP</i> plot and a <i>QQ</i> plot, using the theoretical normal distribution as a reference.
	(cont.) →

Exercise 1.30	(cont.)
	Add Monte Carlo bands from the envelop of 19 simulations. Instead of the uniform distribution, you have to use the normal distribution to generate the Monte Carlo bands. Then you have to represent the results in the co-ordinate system of the <i>QQ</i> plots, that is the x axis represents the quantiles of the normal distribution. <i>Hint:</i> inspect the source of <i>qqnorm()</i> .
*	The bands are initially bands for the standard normal distribution. Find bands adjusted in scale and location of the data at hand.

### 1.5.3 Complements: Functions

R commands can be grouped to functions. Functions may be parametrised by arguments. Functions provide a flexible way of code reusability.

Example of a function:

#### Example 1.14: R Functions

```
ppdemo <- function (x, samps = 19) {  
  # samps: nr of simulations  
  
  y <- (1:length(x))/length(x)  
  plot(sort(x), y, xlab = substitute(x), ylab = expression(F[n]),  
       main = "distribution function with Monte Carlo bands (unif.)",  
       type = "s")  
  mtext(paste(samps, "Monte Carlo samples"), side = 3)  
  samples <- matrix(runif(length(x)* samps),  
                    nrow = length(x), ncol = samps)  
  samples <- apply(samples, 2, sort)  
  envelope <- t(apply(samples, 1, range))  
  lines(envelope[, 1], y, type = "s", col = "red");  
  lines(envelope[, 2], y, type = "s", col = "red")  
}
```

In *ppdemo()* we used the function *mtext()* to generate axis annotations.

Functions are called using the form *<name>(<actual argument list>)*.

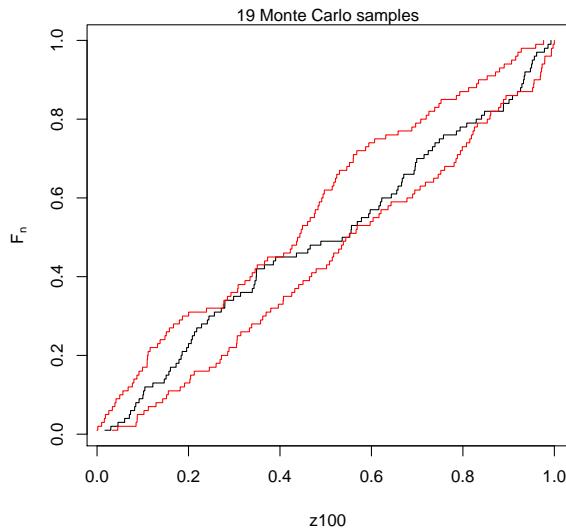
expl:ch01-envfun

expl:ch01-ppdemo

***Example 1.15: R Function Call***

*Input*

```
z100 <- runif(100)
ppdemo(z100)
```

**distribution function with Monte Carlo bands (unif.)**

<b>R Function Declarations</b>	
<b>Declarations</b>	<code>function ( &lt;formal argument list&gt; ) &lt;expression&gt;</code> <i>Example:</i> <code>fak &lt;- function(n) prod(1:n)</code>
<b>Formal argument</b>	<code>&lt;argument name&gt;</code> <code>&lt;argument name&gt; = &lt;default value&gt;</code>
<b>Formal argument list</b>	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 &lt;- function (...)mean(c(...))</code>
<b>Function result</b>	<code>return &lt;value&gt;</code> stops function evaluation and returns value  <code>&lt;value&gt;</code> as last expression in a function declaration: returns value

(cont.)→

<b>R Function Declarations</b> (cont.)	
<b>Function result</b>	<code>&lt;Variable&gt;&lt;-&lt;value&gt;</code> returns value. In general, assignments operate only on local copies of the variable. The assignment with <code>&lt;-</code> however looks for the target in the complete search chain.

<b>R Function Call</b>	
<b>Function call</b>	<code>&lt;name&gt;(&lt;Supplied (actual) argument list&gt;)</code> <i>Example:</i> <code>fak(3)</code>
<b>Supplied argument list</b>	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 for a corresponding actual argument is missing for a formal argument. <i>Syntax:</i> <code>&lt;list of values&gt;</code> <code>&lt;argument name&gt; = &lt;values&gt;</code> <i>Example:</i> <code>rnorm(10, sd = 2)</code>

If only the name of a function is given as input, the definition of the function is returned, that is the function is listed. Example:

pl:ch01-ppdemolist

**Example (cont.): R Function Listing****Example 1.16: R Function Listing***Input*  
ppdemo*Output*

```
function (x, samps = 19) {
  # samps: nr of simulations

  y <- (1:length(x))/length(x)
  plot(sort(x), y, xlab = substitute(x), ylab = expression(F[n]),
       main = "distribution function with Monte Carlo bands (unif.)",
       type = "s")
  mtext(paste(samps, "Monte Carlo samples"), side = 3)
  samples <- matrix(runif(length(x)* samps),
                     nrow = length(x), ncol = samps)
  samples <- apply(samples, 2, sort)
  envelope <- t(apply(samples, 1, range))
  lines(envelope[, 1], y, type = "s", col = "red");
  lines(envelope[, 2], y, type = "s", col = "red")
}
```

a:ch01:e5fun

Exercise 1.31	
	Rework your programming exercises and write reusable parts as functions.

Arguments for functions are passed by value. Each function receives a copy of the actual argument values upon call. This guarantees a safe programming environment. On the other side, this leads to memory requirements and a time penalty. In situations where the argument size is large or time requirements are critical this expenditure can be avoided by direct access to variables which are defined in the function environment. Techniques to achieve this are documented in [R:Gentleman+Thaka:2000](#) [10].

**ToDo:**

example in  
appendix

Functions in R can be nested, that is one “outer” function can contain other “inner” functions. These inner functions are not visible globally, but only in their containing outer function.

Functions can have objects as results. An object is explicitly returned as a result by `return(obj)`. Results can be returned implicitly as well. If the end of a function is reached without calling `return()` the value of the expression evaluated last is returned.

*Input*  
circlearea <- function( r ) r^2 \* pi  
circlearea(1:4)*Output*

```
[1] 3.141593 12.566371 28.274334 50.265482
```

Results can also be provided so that they are only passed upon request. We have met this technique when studying the histogram. The function call `hist(x)` does not return a result, but has the (intended) side effect of drawing a histogram. But if we use `hist()` in an expression, for example in an assignment `xhist <- hist(x)`, we get a description of the histogram returned as the function value. To return results only upon request, instead of `return(obj)` the expression `invisible(obj)` is used.

`a:ch01:fun`

Exercise 1.32	
	<p>Write as functions:</p> <ul style="list-style-type: none"> <li>• A function <code>ehist</code> showing an augmented histogram.</li> <li>• A function <code>eecdfe</code> showing the empirical distribution shows.</li> <li>• A function <code>eqqnorm</code> showing a <i>QQ</i> plot with the standard normal distribution as comparison.</li> <li>• A function <code>eboxplot</code> showing a Box&amp;Whisker-Plot.</li> <li>• A wrapper function <code>eplot</code> showing a plot matrix with these four plots.</li> </ul> <p>Your functions should call the standard functions (or modify them, if necessary) and guarantee that the plots have an adequate complete annotation.</p>

While statements in R are processed stepwise and allow inspection of the results step by step, upon call of a function all statements of the function are executed as a block. This can provide a problem for error diagnosis. R provides possibilities to inspect specific functions, in particular to allow stepwise processing of function. For details see appendix A.13 “Debugging and Profiling” on page A-19.

#### 1.5.4 Complements: Enhancing Graphical Displays

So far the R graphics have been used in rudimentary form only. For serious work the graphics has to be enhanced to make its information identifiable and readable. This implies adding headers, captions, axis labels etc. In R there are “high level” and “low level” graphic functions. “High level” functions generate a new display. Additionally, they provide possibilities to control general graphics parameters for the use in this display.

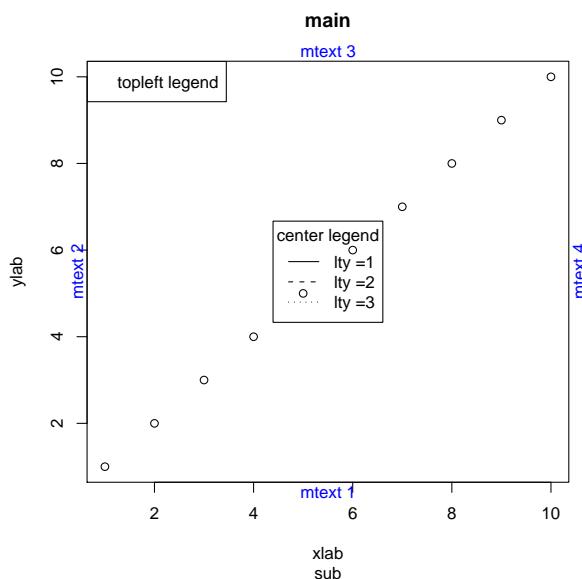
The “low level” functions add elements to an existing plot or modify it in some detail. So for example the function `legend()` can add legends inside a plot.

`expl:ch01-mtext`

**Example 1.17: Margin Text**

*Input*

```
plot(1:10, xlab = "xlab", ylab = "ylab", main = "main", sub = "sub")
mtext("mtext 1", side = 1, col = "blue")
mtext("mtext 2", side = 2, col = "blue")
mtext("mtext 3", side = 3, col = "blue")
mtext("mtext 4", side = 4, col = "blue")
legend("topleft", legend = "topleft legend")
legend("center", legend = c("lty =1", "lty =2", "lty =3"),
       lty = 1:3, title = "center legend")
```



a:ch01:e4plot

Exercise 1.33	
	Use <code>help(plot)</code> to inspect the possibilities to customise the plot function. Some information on detail of the parameters is only available if you use <code>help(plot.default)</code> . Modify your latest plot to have a correct main title.

See also: [R:Rnotes](#) [19] Ch. 12.

For most plot types, annotation, choice of line width and symbols, and colours are at our disposition to enhance the plot. Sometimes, graphic enhancements may be necessary to get a useful plot at all. We give an illustration below, using colour options.

For this example, we need some background knowledge about colour specification. So far, we specified colours by name, that is we use a string such as "`"red"`" from a list of pre-specified

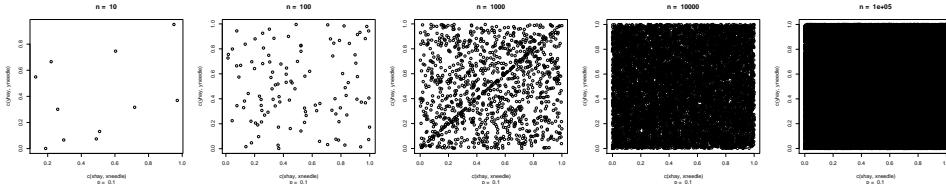
names. The build-in list of colour names is available using `colours()`, but it is also possible to add user defined colour names. The colours themselves are specified by a four byte expression, with the bytes encoding the intensities for the red, green, and blue channel, and the *alpha channel* which is encoding opaqueness. For example, the hexadecimal value `#000000FF` would encode no colour intensity, but full opaqueness. So this would give black. The value `#00000080` would still be black, but only 50% opaque. Drawing once with this colour will appear grey. But if drawings overlay, the colour will add to full black.

Function `col2rgb()` is provided to translate from colour name to a vector of RGB values. Function `rgb()` translates from RGB and opaqueness to colour code.

Now let us look at the example. We construct a needle in the hay stack. The hay stack consists of  $n$  random points in the unit square. The needle is a small number  $p \cdot n$  of random points on the diagonal.

**Example 1.18: Needle in the Hay Stack**

```
NeedleInTheHayStack <- function(n, p=0.1, col="black") {
  xhay <- runif(n); yhay <- runif(n)
  xneedle <- runif(round(p*n)); yneedle <- xneedle
  plot( x = c(xhay, xneedle), y = c(yhay, yneedle), main = paste("n = ", n),
        sub=paste("p = ", p), col= col)
}
oldpar <- par(mfrow=c(1,5))
NeedleInTheHayStack(10); NeedleInTheHayStack(100)
NeedleInTheHayStack(1000); NeedleInTheHayStack(10000)
NeedleInTheHayStack(100000);
par(oldpar)
```



If  $n$  is very small, we have little chance to find the needle in the hay stack. As  $n$  increases, the structure becomes apparent. But if  $n$  gets large, the simple scatter plot seems overloaded and we cannot access the information. But we can enhance the plot using the alpha channel to make it useful again<sup>5</sup>.

**ToDo:**

6 digit numbers

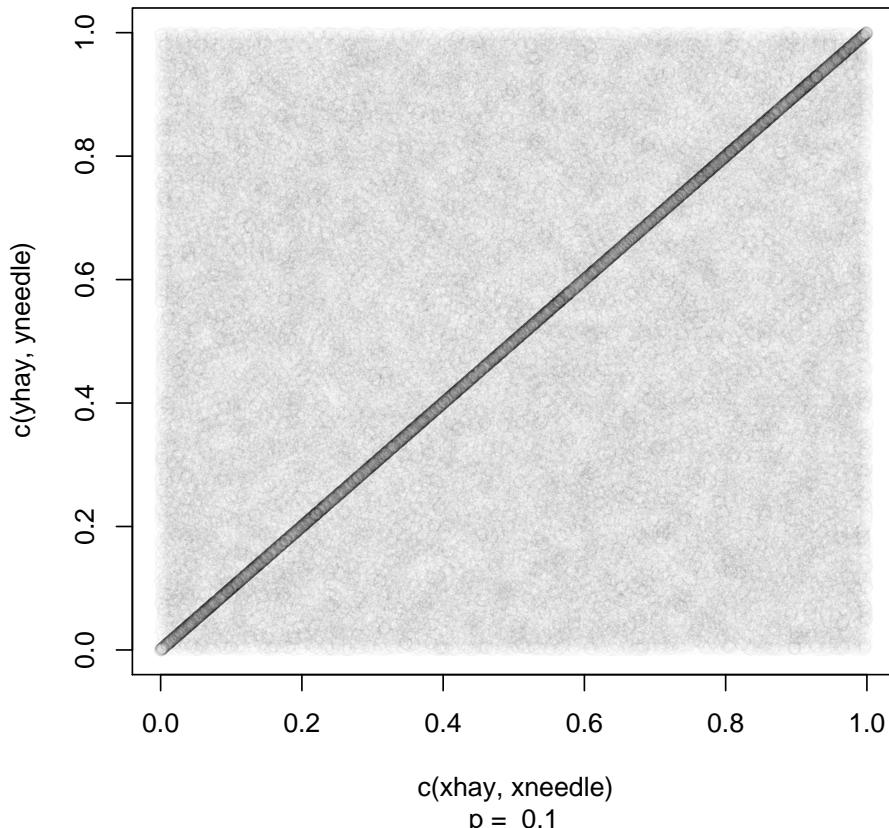
**ToDo:** use TeX image placement

<sup>5</sup> Visualising data sets for large sample sizes is a theme of its own. See [Antony-Unwin2006Graphics-of-LargeDataSets](#).

***Example 1.19: Needle in the Hay Stack***

*Input* \_\_\_\_\_  
`NeedleInTheHayStack(100000, col = rgb(red=0, blue=0, green=0, alpha=0.01))`

**n = 1e+05**



The examples so far used R's basic graphics system. Several other graphics systems are available for the use with R. See for example, the grid/lattice graphics system which we discuss in section 4.1.

### 1.5.5 Complements: R Internals

A typical processing step in R handles a command in three sub-steps

- `parse()` analyses a textual input and translates it into an internal representation as an R expression. R expressions are R objects of a special kind.

- `eval()` interprets this expression and evaluates it. The result is again an R object.
- `print()` displays the resulting object.

Details have to be added:

### *Parse*

The first step consists of two parts, a scanning which scans the input and splits it into tokens, and the parsing proper which takes the tokens and tries to combine them to syntactically correct expressions. The function `parse()` combines both steps. `parse()` can work on local files or data streams as well as on external files denoted by an URL reference.

`parse()` does not resolve names of variables, but just generates an abstract parse tree. Variables in this tree can be resolved using `substitute()`.

As an inverse function `deparse()` is provided. A typical application is to decode the actual arguments of a function call (using `parse()`) and then to derive an informative annotation as in

```
xlab=deparse(substitute(x))
```

### *Eval*

The function `eval()` evaluates an R expression. For this step the references in expressions have to be translated into values, taking into account the actual environment. Since R is an interpreted system, environmental conditions can change. So the evaluation of one expression can lead to different results, depending on the environment.

Each function defines its own local environment. As functions can be nested, so can be the environments. The (nested) environments define a *search path*. Environments can also change as additional libraries are loaded into R. The actual local environment can be queried using `environment()`. With `search()` you get a list of the environments which are searched successively to resolve references. `ls()` gives a list of the objects in an environment.

R's extensibility implies the possibility that names can collide and the translation of names into actual values becomes a problem. As a safeguard, R 2.x provides the possibility to collect names in (guarded) name spaces. In most case this is transparent for the user; the name resolution follows a search path which is determined by the chain of environments. To access an object in a specific name space explicitly, the name space can be specified. So for example `base::pi` can be used as an explicit name for the constant `pi` in the name space `base`), and this access will work even if the user has chosen to introduce some variable with name `pi`, using some arbitrary value.

### *Print*

The function `print()` is implemented as a *polymorphic* function. To execute `print()` R evaluates the class of the objects to print and chooses an appropriate print method. Details are given in section 2.6.5 (page 2-43).

### *Executing Files*

The function `source()` is available to use a file as input for R. The file can be local, or specified by an URL reference. Conventionally for names of R source files, the suffix `.R` is used.

The function `Sweave()` allows to interweave documentation and commands. Conventionally for names of `Sweave()` source files the suffix `.Rnw` is used. Details about the format is in the `Sweave()` documentation

`<http://www.ci.tuwien.ac.at/~leisch/Sweave/Sweave-manual-20060104.pdf>.`

#### *1.5.6 Complements: Packages*

`subs:ch01-lib`

Functions, examples, data sets etc. can be bundled in R packages following certain conventions. The conventions may change with different implementations. The conventions documented in R [23] should be used as a reference. Some packages are part of R by default. Packages for special purposes can be found on the internet, for example using `http://www.cran.r-project.org/src/contrib/PACKAGES.html`.

Packages which are not bundled with R must be installed into the R system first. In general, system specific commands are supplied for this installation. A more convenient way is to do the installation from inside R using `install.packages()`. If no special source is given, `install.packages` tries to access a prepared address (usually the CRAN address given above). But `install.packages` can load packages from any repository specified. In particular, using the form `install.packages(<package>, repos = NULL)`, `<package>` can be an access path for direct access on your local machine.

The function `update.packages()` compares installed versions with the recent state in the repository and updates the version installed locally if necessary.

Once a package is installed, it can be loaded using

`library(pkgname)`

After that, the objects (functions, data sets, ...) defined in the package (functions, data sets, ...) are retrievable on the current search path and can be used directly in any expression.

Packages are released using

`detach(pkgname)`

After this, the package objects do not appear in the current search path any more.

Technically packages are directories, following the R conventions. Usually they are provided in packed and compressed form such as `.tar.gz` files. In the beginning, you will install pre-compiled binary version of packages. If necessary, for example to inspect the source code, the source version can be installed by the same functions.

For the organisation of your own work it is worth following the R conventions and to organise related parts as R packages. R provides several tools supporting package administration if you follow the conventions. The conventions and the tools available to work with packages are documented in [23]. 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`

As a first step: The function `package.skeleton()` helps constructing new packages. `package.skeleton()` does not only generate a package skeleton, but also a help file documenting the next steps to build a loadable package.

Packages must have a file DESCRIPTION with specific information. The Details are documented in [23], and a prototype is generated by `package.skeleton()`. The other components are optional.

Name	Kind	Content
DESCRIPTION	file	A source description following format conventions.
R	directory	R code. Files in this directory should be readable by <code>source()</code> . Recommended name suffix: .R
data	directory	Additional data. Files in this directory should be readable by <code>data()</code> . Recommended name suffixes and formats: .R for R code. Alternatively: .r .tab for tables. Alternatively: .txt, .csv .RData for output of <code>save()</code> . Alternatively: .rda The directory should contain a file <i>00Index</i> with a short survey of the data sets included.
exec	directory	Additional executable files, for example Perl or shell scripts.
inst	directory	Contents is copied (recursively) to the target directory of the installation. In particular, this directory can contain a file CITATION, which is evaluated by the Rfunction <code>citations()</code> .
man	directory	Documentation in the R documentation format (see: [23] "Writing R extensions", available from < <a href="http://www.cran.r-project.org/">http://www.cran.r-project.org/</a> >). Recommended name suffix: .Rd
src	directory	Fortran, C and other sources.
demo	directory	Executable examples. This directory should contain a file <i>00Index</i> with descriptions.

a:ch01:e6lib

Exercise 1.34	
	<p>Install the functions from exercise 1.32 as a package. You can prepare the package with <code>package.skeleton()</code>, if you have already defined the functions.</p> <p>Load the package. Verify that you can still load the package with <code>library()</code> if you have restarted the R system.</p> <p><i>Hint:</i> For an object <math>x</math>, the statement <code>prompt(x)</code> generates a skeleton from which you can build a documentation for <math>x</math>.</p>

## 1.6 Statistical Summary

In this chapter, our leading example was the statistical analysis of a (univariate) random sample. We used a model framework which is central to statistics. The values in the random sample are considered realisations of a random variable, drawn from an underlying theoretical distribution. The aim of the statistical analysis was the inference from the empirical distribution of the random sample on the unknown underlying theoretical distribution. This inference can take two forms: we can compare the empirical distribution with some hypothetical distribution. This is the approach of classical statistics. Or we can attempt to extract features of the underlying distribution from the empirical distribution. This is the data analysis approach.

Both approaches are closely related. The essential tool for both was analysis of the empirical distribution function.

## 1.7 Literature and Additional References

R:Ext

[23]R Development Core Team (2000-2005): Writing R Extensions.

See: <<http://www.r-project.org/manuals.html>>.

gnst77wth

[9] Ganßler, P; Stute, W.: Wahrscheinlichkeitstheorie Heidelberg: Springer 1977.

R:Gentleman+Thaka:2000

[10] Gentleman, R.; Thaka, R.: Lexical Scope and Statistical Computing. Journal of Computational and Graphical Statistics 9 (2000) 491–508.

Generated by Sweave from:

```
Source: /u/math/j40/cvsroot/lectures/src/SIntro/Rintro/Rnw/S01base.Rnw.tex,v
Revision: 1.39
Date: 2008/06/18 16:48:41
name:
Author: j40
$Source: /u/math/j40/cvsroot/lectures/src/SIntro/Rintro/Rnw/S01base.Rnw.tex,v $
$Revision: 1.39 $
$Date: 2008/06/18 16:48:41 $
$name: $
$Author: j40 $
```

---

## CHAPTER 2

---

# Regression

---

ch:02

*Private note:*

*Multiple testing needs to be cleaned up.*

*Point out that usual t-statistics and F-statistics do not account for multiple testing and add proper warnings.*

*Add risk estimation pointing of view and give reference to advanced packages (which?).*

**ToDo:**  
add more  
regression  
graphics

### 2.1 General Regression Model

The (controlled) trial is a common paradigm from experimental sciences: under experimental conditions  $x$  measurement result  $y$  is recorded. This is composed of a systematic effect  $m(x)$  and a measurement error  $\varepsilon$ .

$$y = m(x) + \varepsilon.$$

This is a particular view of the data. Instead of looking at the joint behaviour of  $x$  and  $y$ , an asymmetric view is taken:  $x$  is the “source”,  $y$  (or a change in  $y$ ) the effect. In this view, the “source”  $x$  is considered as given or under direct control;  $y$  is influenced indirectly via the experimental conditions. Experimental conditions should be chosen so the  $\varepsilon$ , the measurement error, is kept as small as possible and vanishes on the average: there should be no systematic error.

From a statistical point of view the main difference of the roles of  $x$  and  $y$  is that for  $y$  the stochastic behaviour is modelled using  $\varepsilon$ , while  $x$  is considered “given” and no stochastics is allotted for  $x$  in the model.

To get a manageable frame, we consider the case where  $x$  can be represented as vector of real variables,  $x \in \mathbb{R}^p$ , and that the measurements gives one dimensional real values,  $y \in \mathbb{R}$ . In a stochastic model we can express the idea formally. One possible formalisation is to represent the measurement error  $\varepsilon$  as a random variable. Assuming additionally that the expected value of  $\varepsilon$  exists, the assumption that the measurement error vanishes on average can be formalised as  $E(\varepsilon) = 0$ .

To analyse the systematic effect  $m$  we consider series of experiments. The index  $i$ ,  $i = 1, \dots, n$ , identifies an experiment in this series, and the model is

$$\begin{aligned} y_i &= m(x_i) + \varepsilon_i \quad i = 1, \dots, n \\ \text{with} \quad x_i &\in \mathbb{R}^p \\ E(\varepsilon_i) &= 0. \end{aligned}$$

Now the statistical problem is: estimate the function  $m$  from the measurement results  $y_i$  at measurement conditions  $x_i$ .

The keywords “curve estimation”, “regression” or “ANOVA” lead to abundant statistical literature about this problem. The following discussion can only be an introduction.

We concentrate mainly on linear regression. This is a very simplified version of the general regression problem. But central aspects of regression can already be illustrated using the linear case. In practical applications, linear models are widely used, and general models are more restricted to sophisticated applications.

To have a unified terminology, we call  $y_i$  the **response** and the components of  $x_{ij}$  with  $j = 1, \dots, p$  are the **regressors**. The function  $m$  is called the **model function**.

Estimators are generally marked by a hat. So  $\hat{m}$  is an estimator for  $m$ .

We introduce two general terms. If we have any estimator, the evaluation at a point  $x$  results in values  $\hat{y} := \hat{m}(x)$ , the **fit** at point  $x$ . If we evaluate an estimator at an observation point, the fit  $\hat{m}(x_i)$  usually will differ from the observation  $y_i$ . The difference

$$R_{x_i}(y_i) := y_i - \hat{m}(x_i)$$

is called **residual**. The residual can be seen as an estimator for the non-observable error term  $\varepsilon$ . Residuals do not exactly match the error terms. This would only be the case if the estimation were exact. The residuals are our only source of information about the error terms. They will be our starting point for inference about the error terms.

The estimation step makes the difference to the situation of independent, identically distributed observations which we discussed in chapter 1. Typically, the estimator  $\hat{m}$  will combine information collected from all experiments in the experimental series. As a consequence, the fitted values, and hence the residuals are random variables which are dependent on all experiments. So even if we could guarantee that the (non-observable) error terms in the original data were independent, the residuals will be dependent variables. We will have to invest additional effort beyond what has been discussed in the previous chapter to compensate for this dependency.

## 2.2 Linear Model

s:2.1

We begin with the regression model - now in vector notation <sup>1</sup> -

$$\begin{aligned} Y &= m(X) + \varepsilon \\ &\text{Y observable random variable with values in } \mathbb{R}^n \\ &X \in \mathbb{R}^{n \times p} \text{ known matrix} \\ &E(\varepsilon) = 0 \text{ non-observable random variable.} \end{aligned} \tag{2.1} \quad \boxed{\text{eq:regrmod}}$$

Additionally we assume that  $m$  is linear. Then there exists a vector  $\beta \in \mathbb{R}^p$  (at least one), so that

$$m(X) = X\beta.$$

<sup>1</sup> We change conventions and notations whenever this change is helpful. The confusion is part of the conventions: in some common conventions, capital letters denote random variable. In other conventions, they denote functions, in still others they denote vectors. We have to leave the resolution to the reader. We try to keep the convention to denote random variables by capital letters, and specific numerical outcomes are denoted by lower case letters. Greek letters are used for parameters.

The regression problem is now reduced to the exercise to estimate  $\beta$  from the information  $(Y, X)$ .

The modified regression model

$$\begin{aligned} Y &= X\beta + \varepsilon \\ Y &\text{ with values in } \mathbb{R}^n \\ X &\in \mathbb{R}^{n \times p} \\ \beta &\in \mathbb{R}^p \text{ unknown vector} \\ E(\varepsilon) &= 0 \end{aligned} \tag{2.2} \quad \boxed{\text{eq:linmod}}$$

is called a *linear model* or *linear regression*. The matrix  $X$  which combines the values of the regressors, hence the information about the experimental conditions, is called the *design matrix* of the model.

**expl:02slr**

**Example 2.1** (Simple Linear Regression) *If the experimental condition is characterised by value of one real variable, which links to the experimental result via a linear model function*

$$m(x) = a + b \cdot x,$$

*we can write a series of experiments with experimental result  $y_i$ , given experimental condition  $x_i$ , in co-ordinates as*

$$y_i = a + b \cdot x_i + \varepsilon_i. \tag{2.3} \quad \boxed{\text{eq:02slrcoord}}$$

*In matrix notation we can write the series of experiments as linear model*

$$Y = \underbrace{\begin{pmatrix} 1 & x_1 \\ \vdots & \vdots \\ 1 & x_n \end{pmatrix}}_X \cdot \underbrace{\begin{pmatrix} a \\ b \end{pmatrix}}_\beta + \varepsilon. \tag{2.4} \quad \boxed{\text{eq:02slrmat}}$$

Simple linear regression is one basic example for linear models. The one way classification is the other basic example.

**expl:oneway**

**Example 2.2** (One Way Classification) *For comparison of  $k$  treatments (in particular for the special case  $k = 2$ ) we use indicator variables which are combined in a matrix. The indicator variable for treatment  $i$  is in column  $i$ . Usually we have repeated observations  $j = 1, \dots, n_i$  under treatment  $i$ , giving a total of  $n = \sum_{i=1}^k n_i$  observations. The model*

$$Y = \underbrace{\begin{pmatrix} 1 & 0 & \dots & 0 \\ \vdots & \vdots & \vdots & \vdots \\ 1 & 0 & \dots & 0 \\ 0 & 1 & \dots & 0 \\ \vdots & \vdots & \vdots & \vdots \\ 0 & 1 & \dots & 0 \\ 0 & 0 & \dots & 1 \\ \vdots & \vdots & \vdots & \vdots \\ 0 & 0 & \dots & 1 \end{pmatrix}}_X \cdot \underbrace{\begin{pmatrix} \mu_1 \\ \mu_2 \\ \vdots \\ \mu_k \end{pmatrix}}_\beta + \varepsilon. \tag{2.5} \quad \boxed{\text{eq:02oneway0mat}}$$

corresponds in co-ordinates

$$y_{ij} = \mu_i + \varepsilon_{ij}. \quad (2.6) \quad \text{eq:02oneway0coord}$$

This is the typical model to test the hypothesis “no difference”  $\mu_1 = \dots = \mu_k$  against the alternative that there is a difference in mean between the treatments.

The same relation can be represented if the measurement values are interpreted as a sum of a basic value  $\mu_0$  and an additional treatment effect  $\mu'_i = \mu_i - \mu_0$ . In co-ordinates, this reads

$$y_{ij} = \mu_0 + \mu'_i + \varepsilon_{ij}. \quad (2.7) \quad \text{eq:02oneway1coord}$$

In matrix notation, this is

$$Y = \underbrace{\begin{pmatrix} 1 & 1 & 0 & \dots & 0 \\ \vdots & \vdots & \vdots & & \vdots \\ 1 & 1 & 0 & \dots & 0 \\ 1 & 0 & 1 & \dots & 0 \\ \vdots & \vdots & \vdots & & \vdots \\ 1 & 0 & 1 & \dots & 0 \\ 1 & 0 & 0 & \dots & 1 \\ \vdots & \vdots & \vdots & & \vdots \\ 1 & 0 & 0 & \dots & 1 \end{pmatrix}}_{X'} \cdot \underbrace{\begin{pmatrix} \mu_0 \\ \mu'_1 \\ \mu'_2 \\ \vdots \\ \mu'_k \end{pmatrix}}_{\beta'} + \varepsilon. \quad (2.8) \quad \text{eq:02oneway1mat}$$

Example 2.2 illustrates that the representation of a problem as a linear model need not be unique. (2.5) and (2.8) are equivalent representations and only the application background can decide which one to prefer.

For mathematical analysis, the design matrix  $X$  is an essential tool. For data analysis, we can use R to generate this matrix (implicitly) for us. R can understand a special notation, the **Wilkinson-Rogers notation**, for model description. With this notation we write

$$y \sim x.$$

The error term is not shown in this representation.

By default, a constant term is assumed. For the one way classification we get model (2.8). If we do not want a constant term (in the regression case for a regression through the origin; in the one way classification the model (2.5) where we do not include an overall mean), we have in co-ordinates

$$y_i = b \cdot x_i + \varepsilon_i.$$

In the Wilkinson-Rogers notation we have to set the constant term to zero explicitly:

$$y \sim 0 + x.$$

Addition regressors can be included using the operator  $+$ . So  $y \sim u + v$  corresponds to the model (in co-ordinates)

$$y_i = a + b \cdot u_i + c \cdot v_i + \varepsilon_i.$$

We return to this notation in section 2.2.4 (page 2-17) and section 2.3 (page 2-20). A summary is given in Appendix A.20 (page A-29).

There is rich literature on linear models. The book “The Theory of Linear Models” of Bent Jørgensen [15] deserves special recommendation. It largely covers the mathematical background of this chapter and contains numerous illustrating examples.

### 2.2.1 Factors

`subs:factor-1`

Using the notation for the description of designs and models, a case oriented description of a design can be machine translated into a design matrix for a linear model in canonical form. Sometimes the translation needs a little bit of coaching. Consider for example a data set

```
y <- c( 1.1, 1.2, 2.4, 2.3, 1.8, 1.9)
x <- c( 1, 1, 2, 2, 3, 3).
```

The vector  $x$  can mean a quantitative vector for the regression model

$$y_i = a + b x_i + \varepsilon_i$$

to be used as a regressor. Or it can be intended for one way classification, the model of one way analysis of variance,

$$y_{ix} = \mu + \alpha_x + \varepsilon_{ix}$$

to be used as indicator for a treatment group. To mark the difference, vectors in R can be defined to be **factors**. Numeric vectors that are no factors are treated as quantitative variable, as in the first example. Factors are treated as indicators and expanded into indicator variables when constructing the design matrix. So

$$y \sim x$$

results in the regression model, but

$$y \sim \text{factor}(x)$$

gives the variance model for the one way layout.

Specifying a parameter `ordered = TRUE` when calling `factor()` marks the variable as ordered. The variable specified will be treated as on an ordinal scale.

$$y \sim \text{factor}(x, \text{ordered} = \text{TRUE})$$

Without this specification, factors are considered to be on a categorial scale.

The values used for factors (both for categorial and ordered) may be mere names and will be substituted internally by numbers. If numerical values are used, they need not be in an ordered sequence.

So

$$\text{factor}( c(2, 2, 5, 5, 4, 4) )$$

results in a vector having three factor values 1, 2, 3, with the names “2”, “5” and “4”. An example using a factor denoted by names would be for example

$$y \sim \text{factor} ( c("Tmt1", "Tmt1", "Tmt2", "Tmt2", "Tmt3", "Tmt3") )$$

The values which can be taken by a factor are called **levels** of the factor. This is an attribute of the factor which can be accessed with `levels()`, as for example in

```
levels(factor( c(2, 2, 5, 5, 4, 4) ))
levels(factor( c("Tmt1", "Tmt1", "Tmt2", "Tmt2", "Tmt3", "Tmt3") ))
```

### 2.2.2 Least Squares Estimation

A first idea for estimation in a linear model can be gained from the following relation: given  $X$ , we have  $E(Y) = X\beta$ . As  $X$  is a matrix, we cannot simply solve this relation for  $\beta$  using a division by  $X$ . But we can expand the relation to  $X^\top E(Y) = X^\top X\beta$ .  $X^\top X$  is a positive symmetric matrix. So at least a pseudo-inverse exists and we can calculate  $(X^\top X)^{-} X^\top E(Y) = \beta$ .

<sup>2</sup> This equation motivates the following estimator:

$$\hat{\beta} = (X^\top X)^{-} X^\top Y. \quad (2.9)$$

**To Do:**

RB:  $\hat{\beta}$  may  
not be  
the only  
LSE or be  
unbiased

Using the model relation  $Y = X\beta + \varepsilon$  <sup>2.2</sup> and  $E(\varepsilon) = 0$ , we get

$$E(\hat{\beta}) = E\left((X^\top X)^{-} X^\top (X\beta + \varepsilon)\right) = \beta, \quad (2.10)$$

eq:02-kqs

eq:02-unbiased

so  $\hat{\beta}$  is an unbiased estimator for  $\beta$ . It is a topic in statistics lectures to discuss whether there are other qualities of this estimator. The Gauss-Markov theorem is a theorem from statistics characterising the estimator <sup>2.2</sup>. We came back to this estimator frequently and give it a name for reference: **Gauss-Markov estimator**. In case of a linear model, such as the regression model, this estimator has a series of optimality properties. For example, this estimator minimises the mean quadratic deviation, that is it is a **least squares estimator** in this model.

The least squares estimator for linear models is implemented as function `lm()`.

We generate an example data set to be used for illustration.

---

*Input*

```
x <- 1:100
err <- rnorm(100, mean = 0, sd = 10)
y <- 2.5*x + err
```

**bsp:02-lm** For this example data set, we get the least squares estimator using

**Example 2.1: Least Squares Estimator**

---

*Input*

```
lm(y ~ x)
```

---

*Output*

```
Call:
lm(formula = y ~ x)

Coefficients:
(Intercept)          x
0.8742        2.4935
```

**a:ch02:lm01**

<sup>2</sup> Here  $X^\top$  means the transposed of the matrix  $X$  and  $(X^\top X)^{-}$  the (Penrose-Moore generalised) inverse of  $(X^\top X)$ .

Exercise 2.1	
	When we generated the data, we did not use a constant term. The model specified for estimation however did not exclude the constant term. Repeat the estimation using the model without a constant term. Compare the results.

The estimator  $\hat{\beta}$  immediately yields an estimation  $\hat{m}$  for the function  $m$  in our original model:

$$\hat{m}(x) = x^\top \cdot \hat{\beta}.$$

The evaluation at the measurement points results in the vector the fitted values  $\hat{Y} = X\hat{\beta}$ .

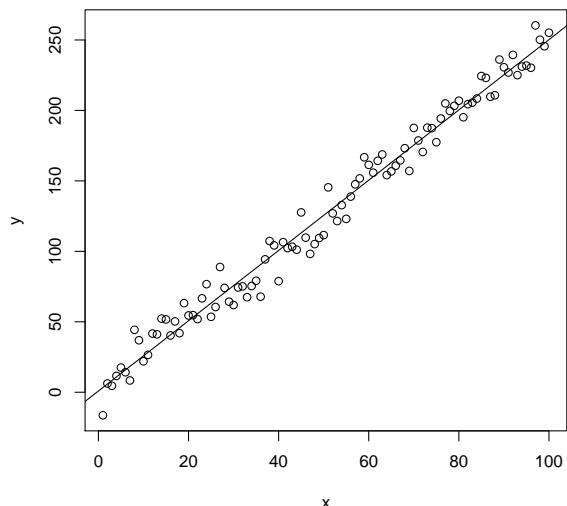
In our example, the fit gives a regression line. Using `plot()` can plot the data points. If we store the result of the regression, we can use it with `abline()` to add the regression line.

expl:02-lmplot

*Example 2.2: Linear Model Plot*

*Input*

```
lmres <- lm(y ~ x)
plot(x, y)
abline(lmres)
```



Function `abline()` is a function to draw lines, using various parametrisations. For more information, see `help(abline)`.

Technically, we can apply the least squares estimation to any data set. The algorithm does not know whether the model assumptions apply, and it does not give us any information about the

quality of the result. It is optimal, but optimality may not mean much if you are in dire straits. To judge the quality of the estimation, we need additional work.

The first step in this direction is to get information about the variance of the estimator. Equation 2.9 tells us that the estimator  $\hat{\beta}$  is a linear function of the observations  $Y$ . The matrix  $X$  is assumed to be known, hence the linear function is considered a known function. So the stochastic variation comes from the error terms contained in  $Y$ , and we have to reconstruct this.

Equation (2.9) tells us how to calculate the fit at the measurement points.

$$\hat{Y} = X(X^\top X)^{-1}X^\top \cdot Y. \quad (2.11) \quad \text{eq:02-fit}$$

The matrix

$$H := X(X^\top X)^{-1}X^\top \quad (2.12) \quad \text{eq:02-hat}$$

is called the **hat matrix**<sup>3</sup>. It is the main tool to analyse the Gauss-Markov estimator for a given design matrix  $X$ . The design matrix, and hence the hat matrix, depends only on the experimental conditions, but not on the result of the experiment. The fit on the other side always refers to a specific outcome of the experiment, the random sample of observed values  $Y$ .

**ToDo:**  
mention  
unbiased

The linear model contains a term  $\varepsilon$ , representing the measurement error or the experimental fluctuation. We cannot observe this error directly. If we could, we would subtract it and get exact information about the model function. But since the error is not observable, we have to resume to indirect inference.

In general, the value of the random observation  $Y$  is different from the fit  $\hat{Y}$ . The difference

$$R_X(Y) := Y - \hat{Y}$$

is called **residual**. The residual can be seen as an estimator for the non-observable error term  $\varepsilon$ . Residuals do not exactly match the error terms. This would only be the case if the estimation were exact. For the general case the relation

$$\begin{aligned} R_X(Y) &= Y - \hat{Y} \\ &= (I - H)Y \\ &= (I - H)(X\beta + \varepsilon) \\ &= (I - H)\varepsilon, \end{aligned} \quad (2.13) \quad \text{eq:02-ewerr}$$

shows that the residuals are linear combinations of the error terms. We have to infer back from these linear combinations of the error term.

If the variance of the error terms does exist, the variance matrix  $\Sigma$  of the error terms  $Var(\varepsilon) = \Sigma$  determines the variance of the residuals:

$$\begin{aligned} Var(R_X(Y)) &= Var((I - H)\varepsilon) \\ &= (I - H)\Sigma(I - H)^\top. \end{aligned} \quad (2.14) \quad \text{eq:02-varerr}$$

So far we only presumed that there is no systematic error. This was formalised as the assumption

$$E(\varepsilon) = 0.$$

We speak of a **simple linear model**, if we have additionally:

$$\begin{aligned} (\varepsilon_i)_{i=1,\dots,n} &\quad \text{are independent} \\ Var(\varepsilon_i) &= \sigma^2 \quad \text{for a } \sigma \text{ not depending on } i. \end{aligned}$$

<sup>3</sup> it puts the hat on top of  $Y$ .

For a linear model we try to estimate the parameter vector  $\beta$ . The variance structure of the vector of error terms introduces nuisance parameters which complicate the estimation. For a simple linear model this nuisance reduces to just one unknown nuisance parameter  $\sigma$ . Equations like [eq:02-varerr](#) 2.14 can be simplified, because now  $\Sigma = \sigma^2 I$  and the parameter  $\sigma$  can be pulled out from the formula. We can estimate this parameter from the residuals, because the *residual variance*

$$s^2 := \frac{1}{n - Rk(X)} \sum_{i=1}^n (Y_i - \hat{Y}_i)^2 \quad (2.15)$$

is an unbiased estimator for  $\sigma^2$ , where  $Rk(X)$  is the rank of the matrix  $X$ . We write  $\widehat{\sigma^2} := s^2$ . (Taking the root is no linear operation and does not preserve the expected value. The residual standard deviation  $\sqrt{s^2}$  is not an unbiased estimator for  $\sigma$ .) Plugged into equation [eq:02-kds](#) 2.9, the residual variance estimator gives an estimator for the variance/covariance matrix of the estimator for  $\beta$ , because in the simple model we have

$$\text{Var}(\hat{\beta}) = \sigma^2 X^\top X \quad (2.16)$$

which can be estimated by  $s^2 X^\top X$ .

If in addition we can assume that the errors have a normal distribution,  $s^2$  and  $\hat{\beta}$  are independent. If we standardise  $\hat{\beta}$  by  $\text{Var}(\hat{\beta})$ , each component has a *t*-distribution, that is we get tests for hypothesis such as  $\beta_j = 0$ .

[bsp:02-lm](#) The standard output in example 2.1 shows only minimal information about the estimator. More information about the estimator, residuals and derived statistics are returned if we ask for a summary.

**bsp:02.summary**

### Example 2.3: Linear Model Summary

<pre>summary(lm( y ~ x))</pre>	<i>Input</i>
<pre>Call: lm(formula = y ~ x)  Residuals:     Min      1Q  Median      3Q     Max  -22.8000 -6.6761 -0.8656  6.6619 23.4974   Coefficients:             Estimate Std. Error t value Pr(&gt; t )     (Intercept)  0.87418   2.04878   0.427    0.67     x           2.49348   0.03522 70.794   &lt;2e-16 ***   --- Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1   Residual standard error: 10.17 on 98 degrees of freedom Multiple R-squared: 0.9808,          Adjusted R-squared: 0.9806  F-statistic: 5012 on 1 and 98 DF,  p-value: &lt; 2.2e-16</pre>	<i>Output</i>

**a:ch02:1**

Exercise 2.2	
	Analyse the output of <code>lm()</code> shown in <code>refReexplbsp:02.summary</code> . Which of the terms can you interpret? Write down your interpretations. For which terms do you need more information?  Generate a commented version of the output.

In section 2.3 (page 2-20) we will present the theoretical background needed to interpret the remaining terms.

A warning needs to be added here. R reports a  $t$ -test value and an error probability for each of the components of the parameter vector  $\beta$ . However the estimation of the components and hence the derived  $t$ -tests are not independent. So to be on the save side, you have to do a **Bonferroni correction**, that is if  $\beta$  has  $p$  components and you want to guarantee an error level of  $\alpha$ , make sure that the nominal levels for your decision are at most  $\alpha/p$ . This is a crude bound to keep you on the save side. In special cases, it may be possible to have finer tools for simultaneous testing. Examples are in section 2.4 (page 2-27).

Calling `lm()` always returns a result if it is appropriate for the data, but it will also return a linear result if the linear model is not adequate. We need additional diagnostics to tell us whether the model is reliable and usable.

Exercise 2.3	
	Let $yy < -2.5 * x + 0.01x^2 + err$ . What are the results you get if you do a regression using the (incorrect) regression model $yy \sim x$ ? Do you get any hints that this model is not adequate?

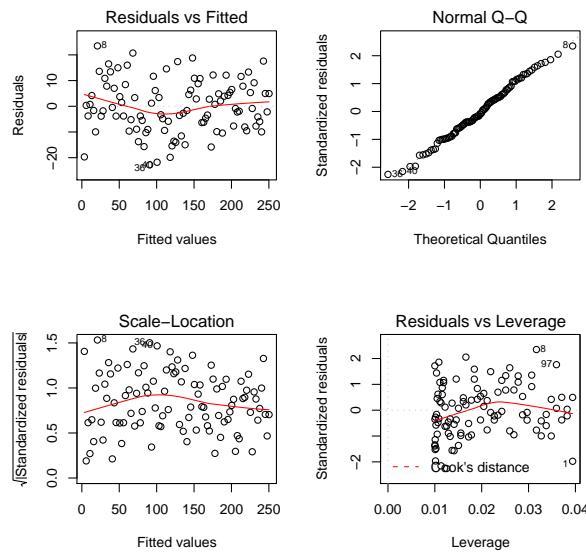
The function `lm()` does not only give an estimation for the linear model, but provides a series of diagnostics that can help to judge whether the model assumptions are acceptable. A representation using `plot()` shows some aspects.

**expl:02-plotlm**

*Example 2.4: Linear Model Plot*

`plot(lm(y ~ x))`

*Input*



The top left plot shows the residuals against the fit. It gives a first survey.

The distribution of the fitted values depends on the design. Unless the design is homogeneous, you cannot expect the residual plot to be homogeneous.

The residuals should approximately look like scatterplot of independent variables. The distribution of the residuals should not vary with the fit. If systematic structures show up in this plot, it is a warning that the model or the model assumptions may not be satisfied.

The previous discussion allows us to be more precise: the residuals should be linear combinations as in (2.13) of independent identically distributed variables. If the model assumptions are satisfied, the variance is given by (2.14).

In a one dimensional situation, a plot of the residuals against the regressor would be sufficient. For  $p$  regressors, the graphical representation becomes difficult. The plot of the residuals against the fit however generalises to higher dimensions of the regressors.

If we start with distribution assumptions about the error terms, we can derive distribution properties of the estimator and of the residuals. The most powerful statements are possible, if the error terms are independent identically distributed with a common normal distribution. In that case, the plot on the upper right should look approximately like a “normal probability plot” of normal random variates, where again “approximately” means: up to transformation with the matrix  $I - H$ .

**a:ch02:1b** The two remaining plots are special diagnostics for linear models (see `help(plot.lm)`).

<b>Exercise 2.4</b>	
	Use <code>plot()</code> to inspect the results of exercise 2.3. Does it give you indications that the linear model is not appropriate? Which indications?

`plot()` provides additional diagnostic plots for linear models. These must be requested explicitly using the parameter `which`.

---

### [help\(lm\)](#)

---

**lm** *Fitting Linear Models*

---

#### *Description*

`lm` is used to fit linear models. It can be used to carry out regression, single stratum analysis of variance and analysis of covariance (although `aov` may provide a more convenient interface for these).

#### *Usage*

```
lm(formula, data, subset, weights, na.action,
   method = "qr", model = TRUE, x = FALSE, y = FALSE, qr = TRUE,
   singular.ok = TRUE, contrasts = NULL, offset, ...)
```

#### *Arguments*

<b>formula</b>	an object of class "formula" (or one that can be coerced to that class): a symbolic description of the model to be fitted. The details of model specification are given under 'Details'.
<b>data</b>	an optional data frame, list or environment (or object coercible by <code>as.data.frame</code> to a data frame) containing the variables in the model. If not found in <code>data</code> , the variables are taken from <code>environment(formula)</code> , typically the environment from which <code>lm</code> is called.
<b>subset</b>	an optional vector specifying a subset of observations to be used in the fitting process.
<b>weights</b>	an optional vector of weights to be used in the fitting process. Should be <code>NULL</code> or a numeric vector. If non- <code>NULL</code> , weighted least squares is used with weights <code>weights</code> (that is, minimizing <code>sum(w*e^2)</code> ); otherwise ordinary least squares is used.
<b>na.action</b>	a function which indicates what should happen when the data contain NAs. The default is set by the <code>na.action</code> setting of <code>options</code> , and is <code>na.fail</code> if that is unset. The 'factory-fresh' default is <code>na.omit</code> . Another possible value is <code>NULL</code> , no action. Value <code>na.exclude</code> can be useful.

<code>method</code>	the method to be used; for fitting, currently only <code>method = "qr"</code> is supported; <code>method = "model.frame"</code> returns the model frame (the same as with <code>model = TRUE</code> , see below).
<code>model, x, y, qr</code>	logicals. If <code>TRUE</code> the corresponding components of the fit (the model frame, the model matrix, the response, the QR decomposition) are returned.
<code>singular.ok</code>	logical. If <code>FALSE</code> (the default in S but not in R) a singular fit is an error.
<code>contrasts</code>	an optional list. See the <code>contrasts.arg</code> of <code>model.matrix.default</code> .
<code>offset</code>	this can be used to specify an <i>a priori</i> known component to be included in the linear predictor during fitting. This should be <code>NULL</code> or a numeric vector of length either one or equal to the number of cases. One or more <code>offset</code> terms can be included in the formula instead or as well, and if both are specified their sum is used. See <code>model.offset</code> .
<code>...</code>	additional arguments to be passed to the low level regression fitting functions (see below).

### Details

Models for `lm` are specified symbolically. A typical model has the form `response ~ terms` where `response` is the (numeric) response vector and `terms` is a series of terms which specifies a linear predictor for `response`. A `terms` specification of the form `first + second` indicates all the terms in `first` together with all the terms in `second` with duplicates removed. A specification of the form `first:second` indicates the set of terms obtained by taking the interactions of all terms in `first` with all terms in `second`. The specification `first*second` indicates the cross of `first` and `second`. This is the same as `first + second + first:second`.

If the formula includes an `offset`, this is evaluated and subtracted from the response.

If `response` is a matrix a linear model is fitted separately by least-squares to each column of the matrix.

See `model.matrix` for some further details. The terms in the formula will be re-ordered so that main effects come first, followed by the interactions, all second-order, all third-order and so on: to avoid this pass a `terms` object as the formula (see `aov` and `demo(glm.vr)` for an example).

A formula has an implied intercept term. To remove this use either `y ~ x - 1` or `y ~ 0 + x`. See `formula` for more details of allowed formulae.

`lm` calls the lower level functions `lm.fit`, etc, see below, for the actual numerical computations. For programming only, you may consider doing likewise.

All of `weights`, `subset` and `offset` are evaluated in the same way as variables in `formula`, that is first in `data` and then in the environment of `formula`.

### Value

`lm` returns an object of class `"lm"` or for multiple responses of class `c("mlm", "lm")`.

The functions `summary` and `anova` are used to obtain and print a summary and analysis of variance table of the results. The generic accessor functions `coefficients`, `effects`, `fitted.values` and `residuals` extract various useful features of the value returned by `lm`.

An object of class `"lm"` is a list containing at least the following components:

<code>coefficients</code>	a named vector of coefficients
<code>residuals</code>	the residuals, that is response minus fitted values.
<code>fitted.values</code>	the fitted mean values.
<code>rank</code>	the numeric rank of the fitted linear model.
<code>weights</code>	(only for weighted fits) the specified weights.
<code>df.residual</code>	the residual degrees of freedom.
<code>call</code>	the matched call.
<code>terms</code>	the <code>terms</code> object used.
<code>contrasts</code>	(only where relevant) the contrasts used.
<code>xlevels</code>	(only where relevant) a record of the levels of the factors used in fitting.
<code>offset</code>	the offset used (missing if none were used).
<code>y</code>	if requested, the response used.
<code>x</code>	if requested, the model matrix used.
<code>model</code>	if requested (the default), the model frame used.

In addition, non-null fits will have components `assign`, `effects` and (unless not requested) `qr` relating to the linear fit, for use by extractor functions such as `summary` and `effects`.

### Using time series

Considerable care is needed when using `lm` with time series.

Unless `na.action = NULL`, the time series attributes are stripped from the variables before the regression is done. (This is necessary as omitting NAs would invalidate the time series attributes, and if NAs are omitted in the middle of the series the result would no longer be a regular time series.)

Even if the time series attributes are retained, they are not used to line up series, so that the time shift of a lagged or differenced regressor would be ignored. It is good practice to prepare a `data` argument by `ts.intersect(..., dframe = TRUE)`, then apply a suitable `na.action` to that data frame and call `lm` with `na.action = NULL` so that residuals and fitted values are time series.

### Note

Offsets specified by `offset` will not be included in predictions by `predict.lm`, whereas those specified by an offset term in the formula will be.

### Author(s)

The design was inspired by the S function of the same name described in Chambers (1992). The implementation of model formula by Ross Ihaka was based on Wilkinson & Rogers (1973).

### References

- Chambers, J. M. (1992) *Linear models*. Chapter 4 of *Statistical Models in S* eds J. M. Chambers and T. J. Hastie, Wadsworth & Brooks/Cole.
- Wilkinson, G. N. and Rogers, C. E. (1973) Symbolic descriptions of factorial models for analysis of variance. *Applied Statistics*, **22**, 392–9.

*See Also*

`summary.lm` for summaries and `anova.lm` for the ANOVA table; `aov` for a different interface.  
 The generic functions `coef`, `effects`, `residuals`, `fitted`, `vcov`.  
`predict.lm` (via `predict`) for prediction, including confidence and prediction intervals;  
`confint` for confidence intervals of *parameters*.  
`lm.influence` for regression diagnostics, and `glm` for **generalized linear models**.  
 The underlying low level functions, `lm.fit` for plain, and `lm.wfit` for weighted regression fitting.  
 More `lm()` examples are available e.g., in `anscombe`, `attitude`, `freeny`, `LifeCycleSavings`, `longley`, `stackloss`, `swiss`.  
`biglm` in package `biglm` for an alternative way to fit linear models to large datasets (especially those with many cases).

*Examples*

```
require(graphics)

## Annette Dobson (1990) "An Introduction to Generalized Linear Models".
## Page 9: Plant Weight Data.
ctl <- c(4.17,5.58,5.18,6.11,4.50,4.61,5.17,4.53,5.33,5.14)
trt <- c(4.81,4.17,4.41,3.59,5.87,3.83,6.03,4.89,4.32,4.69)
group <- gl(2,10,20, labels=c("Ctl","Trt"))
weight <- c(ctl, trt)
anova(lm.D9 <- lm(weight ~ group))
summary(lm.D90 <- lm(weight ~ group - 1))# omitting intercept
summary(resid(lm.D9) - resid(lm.D90)) #- residuals almost identical

opar <- par(mfrow = c(2,2), oma = c(0, 0, 1.1, 0))
plot(lm.D9, las = 1)      # Residuals, Fitted, ...
par(opar)

## model frame :
stopifnot(identical(lm(weight ~ group, method = "model.frame"),
model.frame(lm.D9)))

### less simple examples in "See Also" above
```

To be added to the help information: in the formula notation, with two terms or lists of terms *first* and *second*, *first-second* includes the variables indicated by the first term, but excludes those indicated by second. For more information on the formula notation, see `help(formula)`. A summary is given in appendix A.57 (page A-29).

The hat matrix is a particularity of linear models. Fit and residuals however are general concepts and can be applied for all kind of estimations. Clients are often satisfied seeing a fit (or the estimation). For serious clients, and for statisticians, the residuals often contain more valuable information. They indicate what is not yet covered by the model or the estimation.

### 2.2.3 More Examples for Linear Models

subs:02-BspLinMod

The matrix  $X$  is called the **design matrix** of the model. It can be the matrix of row vectors  $x_i$  representing the measurement conditions for experiment  $i$ . But the design matrix is not restricted to this special case. The seemingly simple class of linear models comprises many import special cases. For example:

*Simple Linear Regression:*

$$y_i = a + b x_i + \varepsilon_i \quad \text{with } x_i \in \mathbb{R}, a, b \in \mathbb{R}$$

can be written as a linear model

$$X = (1 \ x)$$

here  $1 = (1, \dots, 1)^\top \in \mathbb{R}^n$ .

*Polynomial Regression:*

$$y_i = a + b_1 x_i + b_2 x_i^2 + \dots + b_k x_i^k + \varepsilon_i \quad \text{with } x_i \in \mathbb{R}, a, b_j \in \mathbb{R}$$

can be written as a linear model

$$X = (1 \ x \ x^2 \ \dots \ x^k)$$

where  $x^j = (x_1^j \dots x_n^j)^\top$ .

In analogy, a multitude of models can be written as linear models, using other transformations.

*Analysis of Variance: One Way Layout*

Observations are taken under  $m$  experimental conditions, with  $n_j$  observations are taken under condition  $j, j = 1, \dots, m$ . The measurement is the sum of a baseline effect  $\mu$ , a contribution  $\alpha_j$  that is specific for condition  $j$ , and a measurement error

$$y_{ij} = \mu + \alpha_j + \varepsilon_{ij} \quad \text{with } \mu, \alpha_j \in \mathbb{R}, j = 1, \dots, n_i.$$

This can be written as a linear model with  $n = \sum n_j$  and

$$X = (1 \ I_1 \ \dots \ I_m),$$

where  $I_j$  is an indicator variable marking group  $j$ .<sup>4</sup>

*Analysis of Covariance*

Like for the analysis of variance, differences between treatment groups are analysed. But in addition to the treatment, additional influence factors are considered and must be adjusted for. Under experimental condition  $j$ , the observation  $i$  in this group depends on additional influence factors  $x_{ij}$  for observation  $ij$ .

$$y_{ij} = \mu + \alpha_j + b x_{ij} + \varepsilon_{ij} \quad \text{with } \mu, \alpha_j \in \mathbb{R}.$$

<sup>4</sup> For analysis of variance, the convention is to use the first index to mark the group, and the last index to count the observations per group. Index names are given in alphabetical sequence. Following this convention, in comparison to our notation the roles of  $i$  and  $j$  are swapped.

**subs:02-wilkrog**

### 2.2.4 Model Formulae

R allows to specify models by giving the rules for building the design matrix. The syntax for writing these rules is sketched in the description of `lm()`. We will have a closer look at this formula notation. This kind of model specification does not only apply to linear models, but can also be used for more general non-linear models. The model specification is stored as an attribute with attribute name `formula`. It can be manipulated using `formula()`.

*Examples*

$y \sim 1 + x$	corresponds to $y_i = (1 \ x_i)(\beta_1 \ \beta_2)^\top + \varepsilon$ .
$y \sim x$	short for $y \sim 1+x$ (a constant term is assumed implicitly).
$y \sim 0 + x$	corresponds to $y_i = x_i\beta + \varepsilon$ .
$\log(y) \sim x_1 + x_2$	corresponds to $\log(y_i) = (1 \ x_{i1} \ x_{i2})(\beta_1 \ \beta_2 \ \beta_3)^\top + \varepsilon$ (a constant term is assumed implicitly).
<code>lm(y ~ poly(x, 4), data = Experiment)</code>	analyses the data set “Experiment” with a linear model for polynomial regression of degree 4 in $x$ .

There are important special cases for factorial designs:

$y \sim A$	one way analysis of variance with factor $A$ .
$y \sim A + x$	analysis of covariance with factor $A$ and regression covariate $x$ .
$y \sim A + B$	two factor crossed design with factors $A$ and $B$ without interaction.
$y \sim A * B$	two factor crossed design with factors $A$ and $B$ and all interactions (combinations of the levels of $A$ and $B$ ).
$y \sim A/B$	two factor hierarchical layout with factor $A$ and subfactor $B$ .

**a:ch02:2** A table of all operators for model specification formulae is in appendix [A.57](#) (page [A-29](#)).

Exercise 2.5	
	Write the four models from section <a href="#">2.2.3</a> using the R formula notation. <small>subs:02-EspLinMod ta:wrn</small>
	For each of these models, generate an example data set by simulation, and apply <code>lm()</code> to the example. Compare the estimators returned by <code>lm()</code> with the parameters you have used in the simulations.

The model formula used by `lm()` is returned as entry in the function result and recovered from

there. From the formula, R implicitly generates a design matrix. Function `model.matrix()` can be used to inspect design matrix.

**a:ch02:3**

Exercise 2.6	
	<p>Generate three vectors of random variables with a <math>N(\mu_j, 1)</math> distribution, <math>\mu_j = j</math>, <math>j = 1, 3, 9</math>, each of length 10, and combine these into a vector <math>y</math>.</p> <p>Generate a vector <math>x</math> with the values <math>j</math>, <math>j = 1, 3, 9</math>, each repeated 10 times.</p> <p>Calculate the Gauss-Markov estimator in the linear models</p> $y \sim x \text{ and } y \sim \text{factor}(x).$ <p>Inspect the results as a table using <code>summary()</code> and graphically using <code>plot()</code>. Compare the results, and give a written report.</p>

### 2.2.5 Gauss-Markov Estimator and Residuals

We have a closer look on the Gauss-Markov estimator. Knowledge from linear algebra, long thinking or other sources tell us:

**rem:2.1**

#### Remark 2.3

- (1) The design matrix  $X$  defines a mapping  $\mathbb{R}^p \rightarrow \mathbb{R}^n$  with  $\beta \mapsto X\beta$ .  
Let  $\mathcal{M}_X$ ,  $\mathcal{M}_X \subset \mathbb{R}^n$  be the image space of this mapping.  $\mathcal{M}_X$  is the vector space generated by the column vectors from  $X$ .
- (2) If the model assumptions are satisfied,  $E(Y) \in \mathcal{M}_X$ .
- (3)  $\hat{Y} = \pi_{\mathcal{M}_X}(Y)$ , where  $\pi_{\mathcal{M}_X} : \mathbb{R}^n \rightarrow \mathcal{M}_X$  is the (Euclidean) orthogonal projection.
- (4)  $\hat{\beta} = \arg \min_{\beta} |Y - \hat{Y}_{\beta}|^2$  where  $\hat{Y}_{\beta} = X\beta$ .

The characterisation (3) of the Gauss-Markov estimator as an orthogonal projection often helps understanding. The fit is the orthogonal projection of the observation vector on the space of expected values of the model (and hence minimises the quadratic distance). The vector of residuals is the orthogonal complement.

In statistics, the estimator is analysed systematically, and the characterisation given above is just one starting point. Some properties of the estimator can be easily derived using knowledge from probability theory, such as the following lemma:

**th:S02.1**

**Theorem 2.4** Let  $Z$  be a random variable with values in  $\mathbb{R}^n$ , with  $N(0, \sigma^2 I_{n \times n})$  distribution, and let  $\mathbb{R}^n = L_0 \oplus \dots \oplus L_r$  be an orthogonal decomposition. Let  $\pi_i = \pi_{L_i}$  be the orthogonal projection onto  $L_i$ ,  $i = 0, \dots, r$ .  
Then the following holds:

- (i)  $\pi_0(Z), \dots, \pi_r(Z)$  are independent random variables with normal distributions.
- (ii)  $\frac{|\pi_i(Z)|^2}{\sigma^2} \sim \chi^2(\dim L_i)$  for  $i = 0, \dots, r$ .

*Proof.* → probability theory. See for example [Jørgensen 1993, 2.5 Theorem 3].  $\square$

**ToDo:**

std error,  
t-test,  
point wise  
confidence

Using  $\varepsilon = Y - X\beta$  this allows to derive the theoretical distributions for the estimator  $\hat{\beta}$  and the residuals  $Y - \hat{Y}$ .

In particular, for simple linear models, the residual variance can be used to calculate the variance (resp. standard deviation) for each component  $\hat{\beta}_k$ . The corresponding  $t$  statistics and the  $p$ -value for the test of the hypothesis  $\hat{\beta}_k = 0$  are given in the output of `summary()`.

**ToDo:**  
studentised  
residuals,  
pointwise  
confidence  
intervals

Exercise 2.7	
	What is the distribution of $ R_X(Y) ^2 =  Y - \hat{Y} ^2$ , if $\varepsilon$ has a $N(0, \sigma^2 I)$ distribution?

At first sight  $|R_X(Y)|^2 = |Y - \hat{Y}|^2$ , seems an appropriate gauge to judge the quality of a model: small values indicate a good fit, large value indicate a poor fit. However this has to be taken with caution. On one hand this value depends on linear scale factors. On the other hand, the dimension of the spaces involved has to be taken into account.

What happens if additional regressors are taken into the model? We have already seen, that "linear" includes the possibility to model non-linear relations, for example by taking transformed variables into the design matrix. The characterisation (3) in remark [Rem:2.11](#) tells us, that effectively only the vector space spanned by the design matrix is relevant. Here we can see limits for the Gauss-Markov estimator in linear models: if many transformed variables are taken into the model, or generally if the image space determined by the design matrix becomes too large, an over-fitting will result. In the extreme we may get  $\hat{Y} = Y$ . So all residuals are zero, but the estimation is not useful.

**ToDo:**  
zum  
einen..zum  
anderen

We use  $|R_X(Y)|^2 / \dim(L_X)$ , where  $L_X$  is the orthogonal complement of  $\mathcal{M}_X$  in  $\mathbb{R}^n$  (so  $\dim(L_X) = n - \dim(\mathcal{M}_X)$ ) to compensate for the number of dimensions.

a:ch02:5

Exercise 2.8	
	Modify the output of <code>plot.lm()</code> for the linear model so that instead of the Tukey-Anscombe plot the studentised residuals are plotted against the fit.
*	Enhance the <i>QQ</i> -Plot by Monte Carlo bands for independent normal errors. <i>Hint:</i> You cannot generate the bands directly from a normal distribution - you need the distribution of the residuals, not the distribution of the errors.

a:ch02:6

**ToDo:**  
define stu-  
dentised  
residuals

**ToDo:**  
Check  
revise  
Tukey-  
Anscombe

Exercise 2.9	
	<p>Write a procedure that calculates the Gauss-Markov estimator for the simple linear regression</p> $y_i = a + bx_i + \varepsilon_i$ with $x_i \in \mathbb{R}, a, b \in \mathbb{R}$ and shows four plots: <ul style="list-style-type: none"> <li>• response against regressor, with estimated straight line</li> <li>• studentised residuals against Fit</li> <li>• distribution function of the studentised residuals in a <i>QQ</i> plot with confidence bands</li> <li>• histogram of the studentised residuals</li> </ul>

### 2.3 Variance Decomposition and Analysis of Variance

s:2.2

If a simple linear model with normal errors applies, *t*-tests are useful to solve one dimensional problems (tests or confidence intervals for single parameters, point-wise confidence intervals). To solve simultaneous or higher dimensional problems, we need other tools. Instead of differences or mean values which form the basis of *t*-tests, we use norm based distances (as, for example, quadratic distances) which generalise to higher dimensions.

The interpretation of the Gauss-Markov estimator as orthogonal projection (Remark [2.3.3](#)) shows a possibility to compare models: for  $X, X'$  design matrices with  $\mathcal{M}_{X'} \subset \mathcal{M}_X$ , we consider the decomposition  $\mathbb{R}^n = L_0 \oplus \dots \oplus L_r$  with  $L_0 := \mathcal{M}_{X'}$ , and the orthogonal complements  $L_1 := \mathcal{M}_X \ominus \mathcal{M}_{X'}, L_2 := \mathbb{R}^n \ominus \mathcal{M}_X$ . As above  $\pi$  denotes the projection.

$$F := \frac{\frac{1}{\dim(L_1)} |\pi_{\mathcal{M}_X} Y - \pi_{\mathcal{M}_{X'}} Y|^2}{\frac{1}{\dim(L_2)} |Y - \pi_{\mathcal{M}_X} Y|^2}.$$

This statistics, the *F* statistics (in honour of R.A. Fisher) is the basis for the *analysis of variance*, a classical strategy to compare models.

The idea generalises to chains of models. Let  $\mathcal{M}_0 \subset \dots \subset \mathcal{M}_r = \mathbb{R}^n$  be a chain of vector spaces, and  $L_0 := \mathcal{M}_0, L_i := \mathcal{M}_i \ominus \mathcal{M}_{i-1}$  for  $i = 1, \dots, r$  the stepwise orthogonal complements. With the notation as above,

$$\frac{\frac{1}{\dim(L_{i-1})} |\pi_{\mathcal{M}_i} Y - \pi_{\mathcal{M}_{i-1}} Y|^2}{\frac{1}{\dim(L_i)} |Y - \pi_{\mathcal{M}_i} Y|^2}$$

is a test statistics which can be used as a test for the model  $\mathcal{M}_{i-1}$  in comparison to the refined  $\mathcal{M}_i$ .

a:ch02:11

Exercise 2.10	
	What is distribution of <i>F</i> , if $E(Y) \in \mathcal{M}_{X'}$ applies and $\varepsilon$ is distributed as $N(0, \sigma^2 I)$ ?

a:ch02:12

Exercise 2.11	
	<p>Give an explicit formula for the <math>F</math> statistics for analysis of variance in the one way layout</p> $y_{ij} = \mu + \alpha_j + \varepsilon_{ij}$ <p>in comparison to the homogeneous model</p> $y_{ij} = \mu + \varepsilon_{ij}.$

The analysis of variance gives another representation and interpretation of linear models. For example the regression result from [bsp:02.summary](#) gives the following analysis of variance representation

*Example 2.5: Linear Model ANOVA Summary*

<i>Input</i>	<i>Output</i>
<code>summary(aov(lmres))</code>	<pre>Df Sum Sq Mean Sq F value    Pr(&gt;F) x         1 518070  518070  5011.7 &lt; 2.2e-16 *** Residuals 98 10130     103 --- Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1</pre>

[a:ch02:13](#)

Exercise 2.12	
	Analyse the output of <code>lm()</code> shown in <a href="#">refExplbsp:02.summary</a> . Which terms can you interpret now? Give a written report. For which terms do you need more information?

One more item you find in the output is labelled “R-squared”. The term which is given here is an estimator for the fraction of  $Var(Y)$  which is explained by the model:

$$R^2 = \frac{mss}{mss + rss}$$

with  $mss := \frac{1}{n} \sum (\hat{Y}_i - \bar{Y})^2$  and  $rss := \frac{1}{n} \sum (R_X(Y)_i - \bar{R}_X(Y))^2$ . The notation  $R^2$  stems from simple linear regression. In that text, conventionally the correlation  $Cor(X, Y)$  is denoted by  $R$ , and  $R^2 = Cor(X, Y)^2$ .  $R^2$  does not take into account the number of the estimated parameters. As a consequence, it can easily be too optimistic. The term “adjusted R-squared” is re-weighted to take into account the degrees of freedom.

**To Do:**  
**RB:** is  
**X** now  
random?

---

private note >>>

`ans$df <- c(p, rdf, NCOL(Qr$qr))`

```

if (p != attr(z$terms, "intercept")) {
  df.int <- if (attr(z$terms, "intercept"))
    1
  else 0
  ans$r.squared <- mss/(mss + rss)
  ans$adj.r.squared <- 1 - (1 - ans$r.squared) * ((n -
    df.int)/rdf)
  ans$fstatistic <- c(value = (mss/(p - df.int))/resvar,
    numdf = p - df.int, dendf = rdf)
}
else ans$r.squared <- ans$adj.r.squared <- 0
ans$cov.unscaled <- R
dimnames(ans$cov.unscaled) <- dimnames(ans$coefficients)[c(1,
  1)]
if (correlation) {
  ans$correlation <- (R * resvar)/outer(se, se)
  dimnames(ans$correlation) <- dimnames(ans$cov.unscaled)
  ans$symbolic.cor <- symbolic.cor
}

```

---

<<< private note

---



---



---

### **help(anova)**

---

#### **anova**

#### *Anova Tables*

---

#### *Description*

Compute analysis of variance (or deviance) tables for one or more fitted model objects.

#### *Usage*

**anova(object, ...)**

#### *Arguments*

- |               |   |
|---------------|---|
| <b>object</b> | an object containing the results returned by a model fitting function (e.g., <b>lm</b> or <b>glm</b> ). |
| <b>...</b>    | additional objects of the same type.  |

*Value*

This (generic) function returns an object of class `anova`. These objects represent analysis-of-variance and analysis-of-deviance tables. When given a single argument it produces a table which tests whether the model terms are significant.

When given a sequence of objects, `anova` tests the models against one another in the order specified.

The print method for `anova` objects prints tables in a ‘pretty’ form.

*Warning*

The comparison between two or more models will only be valid if they are fitted to the same dataset. This may be a problem if there are missing values and R’s default of `na.action = na.omit` is used.

*References*

Chambers, J. M. and Hastie, T. J. (1992) *Statistical Models in S*, Wadsworth & Brooks/Cole.

*See Also*


---

`coefficients`, `effects`, `fitted.values`, `residuals`, `summary`, `drop1`, `add1`.

---

Models for the analysis of variance can be specified as descriptive rules. The syntax used is the same we have used for regression, that is the Wilkinson-Rogers notation. If terms on the right hand side of the model description are factors, `lm()` will automatically use a variance analysis model instead of a regression model.

The model description determines the linear spaces that contain the expected values for the model and its sub-models. The analysis of variance is not determined uniquely by the model description. The spaces may allow for various orthogonal decompositions (for example, depending on the order in which sub-models are considered). Moreover the specification of the factors for a design determines a system of generators for the spaces. The factors need not be orthogonal. Even independence is not guaranteed.

This holds for all linear models. In regression, dependency is more an exception. With factorial designs, dependency is a frequent phenomenon. One way analysis, written in co-ordinates, may illustrate this problem: with

$$y_{ij} = \mu + \alpha_j + \varepsilon_{ij} \quad \text{with } \mu, \alpha_j \in \mathbb{R}$$

for  $n_j > 0$  the decomposition into  $\mu$  and  $\alpha_j$  is not unique. The underlying reason is: the global factor  $\mu$  defines the vector space spanned by the constant vector 1, and this is contained in the space spanned by the group indicators.

The model formula defines a design matrix  $X$  and hence a model space. An additional matrix  $C$  is used to extend the matrix and to specify a variance decomposition uniquely. The effective design matrix is then  $[1 \ X \ C]$ ;  $C$  is called the **contrast matrix**. The functions for analysis of variance like `lm()` or `aov()` allow for an explicit specification of contrasts.

The function `anova()` operates like a special formatting of the output and is used analogous to `summary()`, for example using the form `anova(lm())`.

a:ch02:14

**ToDo:**  
clean up  
discussion  
of factors

**ToDo:**  
discuss  
choice of  
models  
- model  
dependent  
estimation  
**ToDo:**  
check  
contrasts

Exercise 2.13	
	<p>The file “micronuclei” contains a data set from a mutagenicity test. Cell cultures (50 units each) have been observed in a control group and under 5 chemical treatments. The effect of the substances is to break up the chromosomes and to build micronuclei. Recorded is the size of the micronuclei (relative to the parent nucleus).</p> <p>Read the file “micronuclei” and calculate mean and variance for each group.</p> <p><i>Hint:</i> You can read the file with <code>data()</code>. For files with table format there is special function <code>read.table()</code>. Use <code>help()</code> to get information about both functions.</p> <p>Some selected statistical functions (such as <code>mean</code>) are in table <small>s:A.12 A.22</small> in the appendix.</p>
	<p>Compare the results between the groups. Are treatment effects detectable? <i>Hint:</i> try to formulate the exercise as a one way analysis of variance. You must reshape the data set first, for example with the help of <code>c()</code>.</p>

**ToCh02de:**

micronuclei still available?

**ToDo:** ref  
appendix:  
reshape  
functions  
A:datamamip

a:ch02:16kiwi

Exercise 2.14	
*	Write a function <code>oneway()</code> , which takes a data table as an argument and performs a one way analysis of variance as a test on difference between the columns.
*	Enhance <code>oneway()</code> by adding the necessary diagnostic plots. Which diagnostics are necessary?

Exercise 2.15	
	<p>The industrial enterprise Kiwi Inc.<sup>5</sup> wants to develop a new helicopter for the market. The helicopter design is rated by the time it stays in air before it touches ground<sup>6</sup> from a fixed starting height (ca. 2m). A design drawing is in figure 2.1, page 2-26. What are the factors that can affect the variability of the flight(sink) time? What are the factors that can affect the mean flight duration?</p>
	<p>Perform 30 test flights with a prototype and measure the time in 1/100s. (You will have to co-operate in pairs to carry out the measurements.) Would you consider the recorded times as normally distributed? The requirement is that the mean flight duration reaches at least 2.4s. Does the prototype satisfy the requirement?</p>
	(cont.)→

<sup>5</sup> Following an idea of Alan Lee, Univ. Auckland, New Zealand<sup>6</sup> Kiwis can not fly.

Exercise 2.15	(cont.)
	<p>Your task is to select a design for production. The variants under discussion are:</p> <ul style="list-style-type: none"> <li>rotor width 45mm</li> <li>rotor width 35mm</li> <li>rotor width 45mm with an additional fold for stabilisation</li> <li>rotor width 35mm with an additional fold for stabilisation.</li> </ul> <p>Your budget allows for about 40 test flights. (If you need more test flights, you should give good arguments for this.) Build 4 prototypes, perform the test flights and record the times. Find the design that achieves maximum flight duration. Generate a report. The report should contain the following details:</p> <ul style="list-style-type: none"> <li>• a list of the observed data and a description of the experimental procedure.</li> <li>• suitable plots of the data for each of the designs</li> <li>• an analysis of variance</li> <li>• a clear summary of your conclusions.</li> </ul> <p><i>Additional hints:</i> Randomise the sequence of your experiments. Reduce the variation by providing uniform conditions for the experiment (same height, same launch technique etc.).</p>
	The fold will result in additional production cost. Give an estimate of the gain which can be achieved by this additional cost.

2-26

REGRESSION

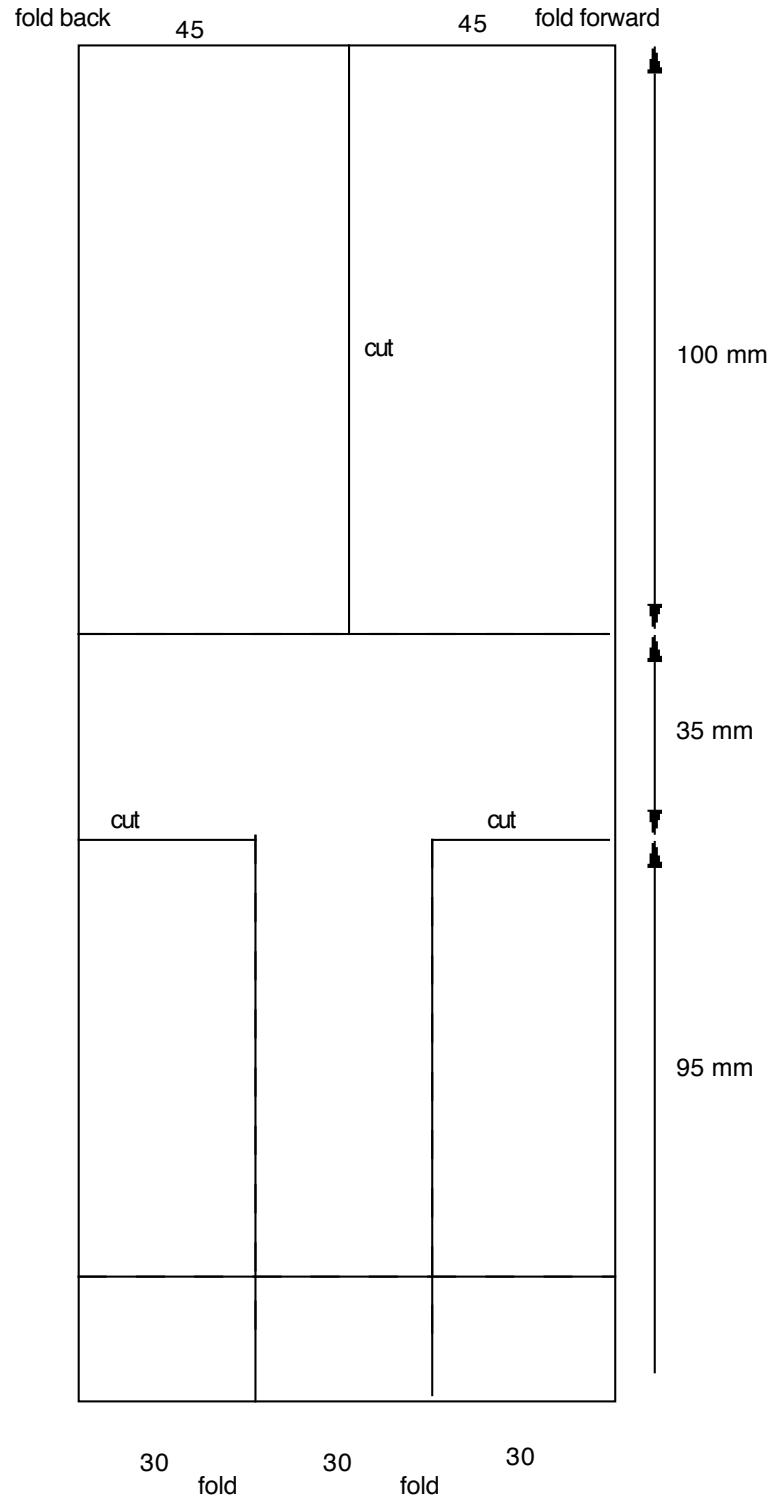


Figure 2.1 *KiwiHopp*

fig:KiwiHopp

## 2.4 Simultaneous Inference

sec:01-simul

### 2.4.1 Scheffé's Confidence Bands

In principle, the least squares estimator estimates all components of the parameter vector simultaneously. On the other hand, the optimality statements in the Gauss-Markov theorem refer only to one dimensional linear statistics. But multivariate confidence statements are possible. The confidence set for confidence level  $1 - \alpha$  derived from the  $F$  distribution has the form

$$\{\hat{\beta} \in \mathbb{R}^k : (\sum_{j=1}^k (\hat{\beta}_j - \beta_j)^2 \|x_j\|^2 / k) / \widehat{\sigma}^2 \leq F_{1-\alpha}(k, n-k)\},$$

that the confidence set is an ellipsoid. Instead of thinking of a quadratic form, we can see the ellipse as the area delimited by all tangents of the ellipsoid. This translates the (single) quadratic condition on the points in the confidence set into (infinitely many) linear conditions. This geometric relation is the core for the following theorem:

**Theorem 2.5** Let  $\mathcal{L} \subset \mathbb{R}^k$  be a linear subspace of dimension  $d$ ;  $EY = Xb$  with  $Rk((X)) = p < n$ . Then

$$P\{\ell^t \beta \in \ell^t \widehat{\beta} \pm (dF_{d,n-\alpha}^\alpha)^{1/2} s(\ell^t (X^t X)^{-1} \ell)^{1/2} \forall \ell \in \mathcal{L}\} = (1 - \alpha).$$

mil8issi  
Proof. [17]2.2, p. 48  $\square$

This is a simultaneous confidence set for all linear combinations from  $\mathcal{L}$ . Translated into a test, this result gives a simultaneous test for all linear hypothesis from  $\mathcal{L}$ . In the special case  $d = 1$  this simultaneous Scheffé test reduces to the usual  $F$ -test. In general, it is not possible to evaluate several tests on the same data material without deflating the confidence level. The  $F$ -test is an exception. After a global  $F$ -test it is possible to test these linear combinations or contrasts individually, without violating the level of the test.

In the case of simple linear regression, we get a confidence ellipsoid in the parameter space, that is the space with co-ordinates *intercept, slope*. If we look at the lines or surfaces defined by parameter values in this confidence ellipsoid, we get a hyperboloid as confidence set for the regression in the regressor/response space.

In the regressor/response space, that is the space of experimental conditions and observations, the interest is often not in a confidence set for the regression, but in a prediction region for observations to come. To get a tolerance area, the stochastic variation by the error term affecting additional observations must be added to the estimation uncertainty. So the tolerance area is increased in comparison to the estimation hyperboloid for the regression. Confidence sets for the regression and tolerance regions for observation can be calculated using function `predict()`. The following figure shows both regions. The function `predict()` is a generic function. For linear models, it calls `predict.lm()`. `predict()` allows to specify new supporting points where the fit is calculated using the estimated model parameters. Variables are matched here by name. `newdata` has to be a data frame `data.frame` with component names corresponding to those of the original variables.

**ToDo:**  
improve  
Scheffe

We prepare a data set as an example.

---

*Input*

```
n <- 100
sigma <- 1
```

```
x <- (1:n)/n-0.5
err <- rnorm(n)
y <- 2.5 * x + sigma*err
lmxy <- lm(y ~ x)
```

**ToDo:**  
check plot

To get better control of the graphics, we calculate the plot limits and supporting points first.

---

*Input*

---

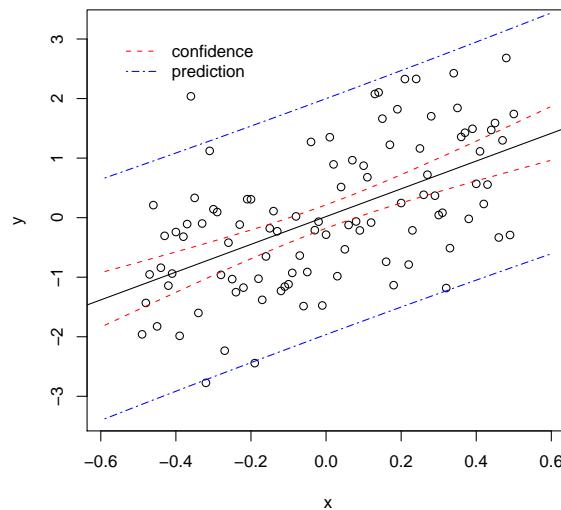
```
plotlim <- function(x){
  xlim <- range(x)
  # check implementation of plot. is this needed?
  del <- xlim[2]-xlim[1]
  if (del>0)
    xlim <- xlim+c(-0.1*del, 0.1*del)
  else xlim <- xlim+c(-0.1, 0.1)
  return(xlim)
}
xlim <- plotlim(x)
ylim <- plotlim(y)
#newx <- data.frame(x = seq(1.5*min(x), 1.5*max(x), 1/(2*n)))
newx <- data.frame(x = seq(xlim[1], xlim[2], 1/(2*n)))
```

For these nodes, we calculate the confidence bands and draw them.

expl:02-lmkonf

**Example 2.6: Linear Model Confidence Bands**

```
plot(x, y, xlim = xlim, ylim = ylim) Input
abline(lmxy)
pred.w.plim <- predict(lmxy, newdata = newx, interval = "prediction")
pred.w.clim <- predict(lmxy, newdata = newx, interval = "confidence")
matplot(newx$x,
        cbind(pred.w.clim[, -1], pred.w.plim[, -1]),
        lty = c(2, 2, 6, 6),
        col = c(2, 2, 4, 4),
        type = "l", add = TRUE)
title(main = "Simultaneous Confidence")
legend("topleft",
       lty = c(2, 6),
       legend = c("confidence", "prediction"),
       col = c(2, 4),
       inset = 0.05, bty = "n")
```

**Simultaneous Confidence****2.4.2 Tukey's Confidence Intervals**

Geometrically the confidence ellipsoid is delimited by its (infinitely many) tangential spaces. Translated into tests this means that infinitely many linear tests are performed simultaneously. In many applications it is possible to formulate more specific questions. For example, to compare  $k$  treatments with treatment effects  $\beta_i$ ,  $i = 1, \dots, k$  we may be interested only in the hypotheses  $\beta_i - \beta_{i'} = 0$ . These reduced question can be formulated in the framework of linear models. When

**ToDo:****Scheffe:**

fix lower/upper plotting bound

**ToDo:**

Comment on Tukey vs Scheffé

specified, they can lead to more powerful tests than the omnibus simultaneous tests. Families of hypotheses of this kind can again be specified as *contrasts*. R supports the use of contrasts for the analysis of variance.

### Case Study: Titre Plates

A typical tool in biology and medicine are titre plates which are used, for example, in cell culture experiments. The plate has small wells in a rectangular grid. Some substrate can be applied to the plate as a whole. Then on each well, a test substance can be applied. Typically, this is done using multi pipettes for a row or a column at once (figure 2.2).

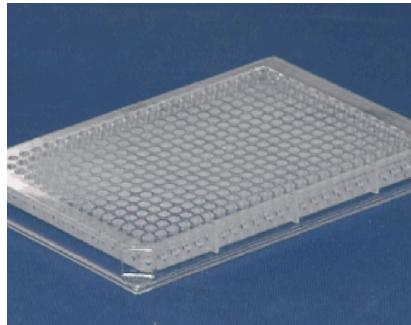


Figure 2.2 *Titre Plates*. With multi pipettes, substances can be applied to a row or a column at once.

`fig:microtitre`

The experiments are often done in series. We will just use one plate from a series as our data example.

---

*Input*

```
p35 <- read.delim("../data/p35.tab")
```

In this experiment, data had been recorded in matrix form. For the analysis with `lm()` we have to transform the data from a matrix form in a long form which has a treatment label in one column. The column *H* of the plate in our experiment contains no treatment, but serves only for comparison between plates in the series for quality control.

---

*Input*

```
s35 <- stack(p35[,3:9])                                # ignore column H
s35 <- data.frame(y=s35$values,
                   Tmt=s35$ind,
                   Lane=rep(1:12, length.out=dim(s35)[1]))    # rename
lmres <- lm(y ~ 0 + Tmt, data= s35)                  # we do not want an overall mean
```

The summary as a linear model gives *t*-tests for the individual coefficients.

---

*Input*

```
summary(lmres)
```

---

Output

```

Call:
lm(formula = y ~ 0 + Tmt, data = s35)

Residuals:
    Min      1Q  Median      3Q     Max 
-0.084833 -0.016354  0.009125  0.022729  0.073083 

Coefficients:
            Estimate Std. Error t value Pr(>|t|)    
TmtA   0.19383   0.01035   18.73   <2e-16 ***
TmtB   0.24892   0.01035   24.06   <2e-16 ***
TmtC   0.23783   0.01035   22.99   <2e-16 ***
TmtD   0.24117   0.01035   23.31   <2e-16 ***
TmtE   0.24392   0.01035   23.57   <2e-16 ***
TmtF   0.23558   0.01035   22.77   <2e-16 ***
TmtG   0.22367   0.01035   21.62   <2e-16 ***  
---
Signif. codes:  0 ‘***’ 0.001 ‘**’ 0.01 ‘*’ 0.05 ‘.’ 0.1 ‘ ’ 1

Residual standard error: 0.03584 on 77 degrees of freedom
Multiple R-squared:  0.9787,    Adjusted R-squared:  0.9768 
F-statistic: 506.2 on 7 and 77 DF,  p-value: < 2.2e-16

```

For this example, tests for individual coefficients are not adequate. `anova()` gives a summary which is tailored to the needs of an analysis of variance.

---

*Input*

---

`anova(lmres)`

---

Output

```

Analysis of Variance Table

Response: y
          Df Sum Sq Mean Sq F value    Pr(>F)    
Tmt       7 4.5513  0.6502  506.15 < 2.2e-16 ***
Residuals 77 0.0989  0.0013
---
Signif. codes:  0 ‘***’ 0.001 ‘**’ 0.01 ‘*’ 0.05 ‘.’ 0.1 ‘ ’ 1

```

Provided the assumptions of the Gauss-linear model are satisfied, this summary says that the treatment effect is significant. The next question that arises immediately is which of the treatments differ significantly, that is we are interested in the contrasts that describe the differences between treatments. Without violating the level of the test, we can use Tukey's approach to post-hoc tests. This is supported by function `glht()` for tests of generalised linear hypotheses, which is provided in the package `multcomp` for multiple testing.

expl:02-lhtres

**Example 2.7: Tukey's Multiple Comparison**

Input		Output	
<code>library(multcomp)</code>			
<code>lhtres&lt;-glht(lmres,linfct=mcp(Tmt="Tukey"))</code>			
<code>summary(lhtres) # multiple tests</code>			
		Simultaneous Tests for General Linear Hypotheses	
Multiple Comparisons of Means: Tukey Contrasts			
Fit: lm(formula = y ~ 0 + Tmt, data = s35)			
Linear Hypotheses:			
	Estimate	Std. Error	t value p value
B - A == 0	0.055083	0.014632	3.765 0.00572 **
C - A == 0	0.044000	0.014632	3.007 0.05252 .
D - A == 0	0.047333	0.014632	3.235 0.02892 *
E - A == 0	0.050083	0.014632	3.423 0.01643 *
F - A == 0	0.041750	0.014632	2.853 0.07806 .
G - A == 0	0.029833	0.014632	2.039 0.39943
C - B == 0	-0.011083	0.014632	-0.757 0.98819
D - B == 0	-0.007750	0.014632	-0.530 0.99832
E - B == 0	-0.005000	0.014632	-0.342 0.99986
F - B == 0	-0.013333	0.014632	-0.911 0.96973
G - B == 0	-0.025250	0.014632	-1.726 0.60107
D - C == 0	0.003333	0.014632	0.228 0.99999
E - C == 0	0.006083	0.014632	0.416 0.99958
F - C == 0	-0.002250	0.014632	-0.154 1.00000
G - C == 0	-0.014167	0.014632	-0.968 0.95930
E - D == 0	0.002750	0.014632	0.188 1.00000
F - D == 0	-0.005583	0.014632	-0.382 0.99974
G - D == 0	-0.017500	0.014632	-1.196 0.89358
F - E == 0	-0.008333	0.014632	-0.570 0.99748
G - E == 0	-0.020250	0.014632	-1.384 0.80863
G - F == 0	-0.011917	0.014632	-0.814 0.98279
---			
Signif. codes: 0 ‘***’ 0.001 ‘**’ 0.01 ‘*’ 0.05 ‘.’ 0.1 ‘ ’ 1			
(Adjusted p values reported -- single-step method)			

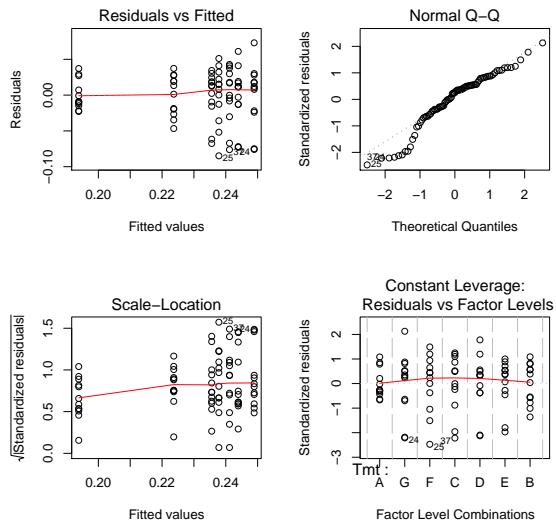
Under the assumptions of the model the significance of the differences between  $A$  and  $B$ , the differences between  $A$  and  $D$  are established.

To check the assumptions, we can use the residuals which can be inspected with `plot()`.

expl:02-plotlhtres

*Input*

```
oldpar <- par(mfrow=c(2,2))
plot(lmres)
par(oldpar)
```



The distribution shows a clear deviation from the normal distribution, in particular at low values. We can inspect these selectively.

expl:02-s35studres

**Example 2.9: Studentized Residuals**

```
#diagnostic
library(MASS)
s35$studres <- studres(lmres)
s35[s35$studres < -1,]
```

	y	Tmt	Lane	studres
13	0.174	B	1	-2.239390
24	0.173	B	12	-2.271296
25	0.153	C	1	-2.559766
33	0.202	C	9	-1.044865
36	0.186	C	12	-1.523409
37	0.165	D	1	-2.279286
48	0.174	D	12	-1.994858
49	0.171	E	1	-2.175828
60	0.172	E	12	-2.144169
61	0.168	F	1	-2.007887
72	0.174	F	12	-1.821449
73	0.177	G	1	-1.367608
84	0.189	G	12	-1.010381

Input

Output

The pattern is conspicuous. Nearly all small values occur at the boundary of the plate.

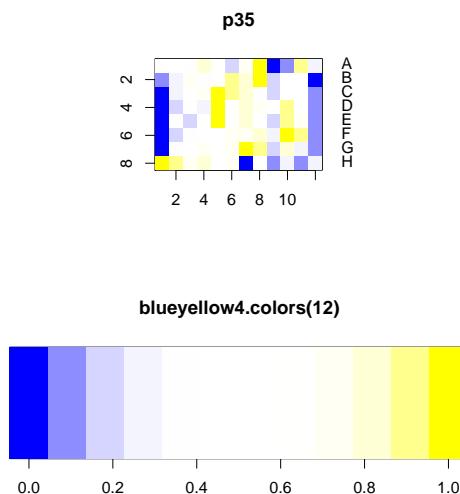
We could as well have detected this pattern by visual inspection:

expl:02-bertin

*# visualisation*  
*# image, the easy way*

```
# image, the easy way
a35 <- as.matrix(p35[3:10])
a35rk <- apply(a35, 2, rank)
#image(a35rk)

# enhanced image, using bertin
library(bertin)
oldpar<- par(mfrow=c(2,1))
image(t(a35rk), col=blueyellow4.colors(12), main="p35")
colramp(blueyellow4.colors(12),12, horizontal=TRUE)
par(oldpar)
```



For independent errors we would expect a random distribution of the range values over the rows. The concentration of the extreme values in the extreme columns shows an inhomogeneity in the experimental production process.

So, for this example, we can report that apparently there is a difference between treatment A and particular other treatments. The finding has to be seen with precaution: the model assumptions are not satisfied. There is a noticeable inhomogeneity between the rows. More important than the findings on the treatment differences may be the hint to look into the production process for the sources of this inhomogeneity.

**ToDo:**  
 rows,  
 columns,  
 lanes?

We would not be able to detect a similar effect for the rows. Since the treatments are applied row-wise, treatment effect and column effect are confounded, and the effects are not identifiable.

## 2.5 Beyond Linear Regression

s:2.3

### Transformations

subs:3.4.1

Finding scale and location parameters can be understood as attempt to transform the distribution to some reference form. However scale and location cover only linear transformations.

The **Box -Cox transformations**

$$y^{(\lambda)} = \begin{cases} \frac{y^\lambda - 1}{\lambda} & \text{for } \lambda \neq 0, \\ \log(y) & \text{for } \lambda = 0 \end{cases}$$

are a family which is scaled to imbed the logarithm transformation smoothly in the power transformations. The function `boxcox()` in `library(MASS)` can be used to choose  $\lambda$ .

Using transformations, linear models can sometimes be used to cover non-linear situations.

Generalised linear models are an extension of linear models that allow to cover certain transformations within the model. For more information, see [31].

#### 2.5.1 Generalized 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 count among the statistical models which are best investigated. Theory and algorithms are far advanced. So it is inviting to try how this class of models can be extended while still allowing to use the theoretical and algorithmical know-how.

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

### 2.5.2 Local Regression

We now make a big jump. We have discussed linear models. We know that we can represent non-linear functions using linear models, but the term entering the functions must be specified off hand. Too many terms lead to over-fitting; so it is not the best idea to leave it to the algorithm for linear models to do the model selection. Statistical treatment of a regression problem with mild assumptions on the model function remains a problem.

A partial approach for a solution comes from analysis. In analysis, it is a standard technique to use local approximations for functions. The analogous approach is to use a localised variant instead of a global estimation. We still assume that

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

but assume linearity only locally:

$$m(x) \approx x' \beta_{x_0} \quad \beta_{x_0} \in \mathbb{R}^p \text{ and } x \approx x_0.$$

For practical work, abstract asymptotics is not sufficient. We have to specify what  $\approx$  means. This can be done with reference to some scales (for example  $x \approx x_0$  if  $|x - x_0| < 3$ ) or with reference to the design (for example  $x \approx x_0$  if  $\#\{i : |x - x_i| \leq |x - x_0| < n/3\}$ ). Up to date implementations use refinements that cannot be discussed here. For an illustration of concepts, consider:

Localised Gauss-Markov estimator:

For  $x \in \mathbb{R}^p$ , find

$$\delta = \min_d : (\#\{i : |x - x_i| \leq d\}) \geq n \cdot f$$

where  $f$  is a chosen fraction (for example 0.5).

Calculate the Gauss-Markov estimator  $\hat{\beta}_x$ , using only those observations for which  $|x - x_i| \leq \delta$ . Estimate

$$\hat{m}(x) = x' \hat{\beta}_x.$$

This coarsening ignores all measurement points which have a distance of more than  $\delta$ . More subtle methods use a weight function to reduce the influence of distant measurement points progressively.

[help\(loess\)](#)

`loess`

*Local Polynomial Regression Fitting*

#### Description

Fit a polynomial surface determined by one or more numerical predictors, using local fitting.

*Usage*

```
loess(formula, data, weights, subset, na.action, model = FALSE,
      span = 0.75, enp.target, degree = 2,
      parametric = FALSE, drop.square = FALSE, normalize = TRUE,
      family = c("gaussian", "symmetric"),
      method = c("loess", "model.frame"),
      control = loess.control(...), ...)
```

*Arguments*

<b>formula</b>	a formula specifying the numeric response and one to four numeric predictors (best specified via an interaction, but can also be specified additively). Will be coerced to a formula if necessary.
<b>data</b>	an optional data frame, list or environment (or object coercible by <code>as.data.frame</code> to a data frame) containing the variables in the model. If not found in <b>data</b> , the variables are taken from <code>environment(formula)</code> , typically the environment from which <b>loess</b> is called.
<b>weights</b>	optional weights for each case.
<b>subset</b>	an optional specification of a subset of the data to be used.
<b>na.action</b>	the action to be taken with missing values in the response or predictors. The default is given by <code>getOption("na.action")</code> .
<b>model</b>	should the model frame be returned?
<b>span</b>	the parameter $\alpha$ which controls the degree of smoothing.
<b>enp.target</b>	an alternative way to specify <b>span</b> , as the approximate equivalent number of parameters to be used.
<b>degree</b>	the degree of the polynomials to be used, up to 2.
<b>parametric</b>	should any terms be fitted globally rather than locally? Terms can be specified by name, number or as a logical vector of the same length as the number of predictors.
<b>drop.square</b>	for fits with more than one predictor and <b>degree=2</b> , should the quadratic term (and cross-terms) be dropped for particular predictors? Terms are specified in the same way as for <b>parametric</b> .
<b>normalize</b>	should the predictors be normalized to a common scale if there is more than one? The normalization used is to set the 10% trimmed standard deviation to one. Set to false for spatial coordinate predictors and others know to be a common scale.
<b>family</b>	if "gaussian" fitting is by least-squares, and if "symmetric" a re-descending M estimator is used with Tukey's biweight function.
<b>method</b>	fit the model or just extract the model frame.
<b>control</b>	control parameters: see <code>loess.control</code> .
<b>...</b>	control parameters can also be supplied directly.

*Details*

Fitting is done locally. That is, for the fit at point  $x$ , the fit is made using points in a neighbourhood of  $x$ , weighted by their distance from  $x$  (with differences in ‘parametric’ variables being ignored when computing the distance). The size of the neighbourhood is controlled by  $\alpha$  (set by `span` or `enp.target`). For  $\alpha < 1$ , the neighbourhood includes proportion  $\alpha$  of the points, and these have tricubic weighting (proportional to  $(1 - (\text{dist}/\text{maxdist})^3)^3$ ). For  $\alpha > 1$ , all points are used, with the ‘maximum distance’ assumed to be  $\alpha^{1/p}$  times the actual maximum distance for  $p$  explanatory variables.

For the default family, fitting is by (weighted) least squares. For `family="symmetric"` a few iterations of an M-estimation procedure with Tukey’s biweight are used. Be aware that as the initial value is the least-squares fit, this need not be a very resistant fit.

It can be important to tune the control list to achieve acceptable speed. See `loess.control` for details.

*Value*

An object of class “`loess`”.

*Note*

As this is based on the `cloess` package available at `netlib`, it is similar to but not identical to the `loess` function of S. In particular, conditioning is not implemented.

The memory usage of this implementation of `loess` is roughly quadratic in the number of points, with 1000 points taking about 10Mb.

*Author(s)*

B. D. Ripley, based on the `cloess` package of Cleveland, Grosse and Shyu available at <http://www.netlib.org/a/>.

*References*

W. S. Cleveland, E. Grosse and W. M. Shyu (1992) Local regression models. Chapter 8 of *Statistical Models in S* eds J.M. Chambers and T.J. Hastie, Wadsworth & Brooks/Cole.

*See Also*

`loess.control`, `predict.loess`.  
`lowess`, the ancestor of `loess` (with different defaults!).

*Examples*

```
cars.lo <- loess(dist ~ speed, cars)
predict(cars.lo, data.frame(speed = seq(5, 30, 1)), se = TRUE)
# to allow extrapolation
cars.lo2 <- loess(dist ~ speed, cars,
control = loess.control(surface = "direct"))
predict(cars.lo2, data.frame(speed = seq(5, 30, 1)), se = TRUE)
```

While linear regression is forced by the model assumptions to give a linear image (or an image with linear parametrisation) under all circumstances, a localised variant can represent non linearities properly. The analysis of this family of methods gives a branch of its own in statistics: nonparametric regression.

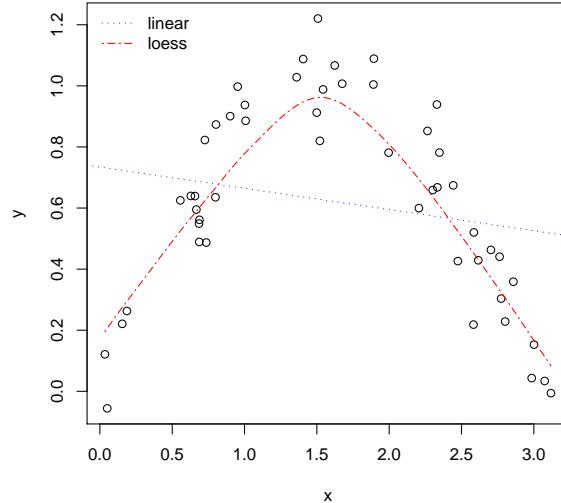
We prepare an example:

```
x <- runif(50) * pi
y <- sin(x)+rnorm(50)/10
```

expl:02-loess

#### Example 2.11: Non-Linear Regression

```
plot(x, y)
abline(lm(y ~ x), lty = 3, col = "blue")
lines(loess.smooth(x, y), lty = 6, col = "red")
legend("topleft",
       legend = c("linear", "loess"),
       lty = c(3, 6), col = c("blue", "red"), bty = "n")
```



## 2.6 R Complements

s:2.4

### 2.6.1 Complements: Discretisation

In analogy to the construction of histograms we can discretise the data. We did this already with respect to the regressors when we discussed the helicopter example. The response can be discretised as well. This converts a regression problem into a contingency table problem. We will not delve into this possibility here.

**ToDo:**  
ch02: Add example

### 2.6.2 Complements: External Data

Data, as well as any other R objects can be written into an external file using `save()`, and read from there again using `load()`. Data are stored in these files in compressed form. The R storage files are exchangeable between different R systems.

Data In/Output for R	
<code>save()</code>	stores data in an external file. Syntax: <code>save(&lt;names of the objects to store&gt;, file = &lt;file name&gt;, ...)</code>
<code>load()</code>	loads data from an external file. Syntax: <code>load(file = &lt;file name&gt;, ...)</code>

Often data are collected or prepared with other systems. R provides a series of functions to read data in various formats. See Appendix A.17 (page A-25). More information is given in the manual “Data Import/Export” ([20]).

The function `data()` bundles various access routines, provided access paths and file names follow the R conventions.

Often imported data have to be post-processed to align the format to the calling conventions of the R functions which will be applied.

For example, `lm` expects the regressors as separate variables. On the other side, for factorial designs it is usual to collect the results in a table where the factor-levels appear as row or column labels. The function `stack()` transforms tables into columns. See also Appendix A.9 (page A-12).

### 2.6.3 Complement: Testing Software

All algorithms, like the algorithms for linear models and their variants, should be treated with the same precautions as any mathematical publication or quotation. Unfortunately even simple programs soon have a semantic complexity far beyond that of many mathematical proof. The usual strategies like re-computing or stepwise execution forbid themselves. Selective testing

must take the place of complete examination. An example of how to design test strategies is documented in [26].

a:ch02:211 Scrutiny is necessary for the implementation as well as for the underlying abstract algorithm.

Exercise 2.16	
	For this series of exercises, let $y_i = a + bx_i + \varepsilon_i$ with $\varepsilon_i$ iid $\sim N(0, \sigma^2)$ and $x_i = i, i = 1, \dots, 10$ .
	Chose a strategy to inspect <code>lm()</code> with regard to the parameter space $(a, b, \sigma^2)$ . Are there apparent cellular decompositions for the parameters $a, b, \sigma^2$ ? What are the trivial cases? What are the asymptotics that apply? Choose test points in the interior of each cell and on the boundaries. Perform these tests and summarise the results.
	What are the symmetries/anti-symmetries that apply? Check for these symmetries.
	Which invariant or covariate behaviour applies? Check for these invariant or covariate behaviour.

a:ch02:221

Exercise 2.17	
	For this series of exercises, let $y_i = a + bx_i + \varepsilon_i$ with $\varepsilon_i$ iid $\sim N(0, \sigma^2)$ .
	What are the extremal designs ( $x_i$ )? Check the behaviour of <code>lm()</code> for four extremal designs.
	Perform the tests from the last exercise, now with variable design. Summarise your results.

a:ch02:23

Exercise 2.18	
	For this series of exercises, let $y_i = a + bx_i + \varepsilon_i$ with $\varepsilon_i$ iid $\sim N(0, \sigma^2)$ .
	Modify <code>lm()</code> to give a fail-safe function for simple linear models which checks deviations from the model assumptions as well.

## 2.6.4 R Data Types

R is an interpreted programming language. It tries to allow a flexible handling of definitions and specifications for the user. For performance reasons, R attempts to evaluate expressions and terms as late as possible. This requires some restrictions for the language which makes R differ from other programming languages.

**ToDo:**  
conditional,  
function  
parameters,  
promises

R does not have abstract data types. A data type is defined by its instances, the variables.

The data type of a variable is dynamic: in the same context the same name can refer to different variable values and different variable types at different times.

At any time however any variable has a well defined type. The R type system is best understood with reference to its historical development and the corresponding functions. Initially the type was described by `mode()` (for example “numeric”) and `storage.mode()` (for example “integer” or “real”).

Essentially both functions have been replaced by `typeof()`. A summary of the information reported by `typeof()` is in [21].

More complex data types are derived from the basic types documented in [21] by attaching attributes to the variables. This is done using function `attr()` which can also be applied to inspect attributes. So for example a matrix or an array are just special vectors, which are special by having a `dim` attribute. The `class` attribute served to denote and name a class explicitly.

Type test and conversion functions are supplied for the main types: `is.<type>()` tests for a type, `as.<type>()` converts to a type.

For the exceptional values, there are special test functions `is.inf()`, `is.na()` and `is.nan()`. A `NaN` value is always considered as missing, so both `is.na(NaN)` and `is.nan(NaN)` will return `TRUE`.

See also Appendix A.4 (page A-4)..

### 2.6.5 Classes and Polymorphic Functions

ch02:polymorph

As the system developed on, concepts from object oriented programming have been taken into R. The special attribute `class` is used here. The name of the type (or the “class”) is stored as `class` attribute. Multiple class membership is possible. In this case, `class` is a vector of class names. So for example an ordered factor has the attribute `class = c("ordered", "factor")`. For class management, the functions `class()`, `unclass()`, `inherits()` are provided.

Class membership in R is based on trust. R does not check, whether a data structures is consistent with its acclaimed class.

Functions like `plot()`, `print()` and many other inspect the class or type of its argument and branch to specialised functions as appropriate. This is called function polymorphism. If you list a polymorphic function, first you get only a clue the apart from special cases a dispatch function `UseMethod()` is called. Example:

expl:02-plot

**Example 2.12: Method Dispatch**

<i>Input</i>	<i>Output</i>
<code>plot</code>	
	<code>function (x, y, ...)</code>
	{
	<code>if (is.function(x) &amp;&amp; is.null(attr(x, "class")))</code>
	<code>if (missing(y))</code>
	<code>y &lt;- NULL</code>
	<code>hasylab &lt;- function(...) !all(is.na(pmatch(names(list(...)),</code>
	<code>"ylab"))</code>
	<code>if (hasylab(...))</code>
	<code>plot.function(x, y, ...)</code>
	<code>else plot.function(x, y, ylab = paste(deparse(substitute(x)),</code>
	<code>"(x)" ), ...)</code>
	<code>}</code>
	<code>else UseMethod("plot")</code>
	}
	<code>&lt;environment: namespace:graphics&gt;</code>

`UseMethod()` tries to identify the class of the first argument that was used in the function call. Then it tries to find a specialisation of the function, a method, for this class and finally it calls this function. For **polymorphic** functions the available methods are listed with `methods()`, for example `methods(plot)`.

### 2.6.6 Extractor Functions

Functions like `lm()` yield complex data with rich information. In a purely object oriented environment, access methods and data would be encapsulated. In R, object orientation is implemented rudimentarily and with varying models. This partially reflects the historical development of the language. For sufficiently general structures, access methods as in section 2.6.5 (page 2-43) are available. For the objects as returned by `lm()`, for example, a series of extractor functions all provide access to components in various forms.

<i>Extractor Functions for lm</i>	
<code>coef()</code>	extracts the estimated coefficients.
<code>effects()</code>	extracts successive orthogonal components.
<code>residuals()</code>	raw residuals.
<code>stdres()</code>	(in <code>library(MASS)</code> ) standardised residuals.

(cont.)→

<i>Extractor Functions for lm</i> (cont.)	
<i>studres()</i>	(in <i>library(MASS)</i> ) externally standardised residuals.
<i>fitted()</i>	
<i>vcov()</i>	variance/covariance matrix of the estimated parameter.
<i>predict()</i>	confidence and tolerance intervals.
<i>confint()</i>	confidence intervals for parameter.
<i>influence()</i>	extracts influence diagnostics.
<i>model.matrix()</i>	derives the design matrix.

## 2.7 Statistical Summary

The leading example in this chapter was the statistical analysis of a functional relationship. The models considered are finite in the sense that only a finite dimensional function space was taken into account to describe the relation between regressors and response. The stochastic component in these models was restricted to a (one dimensional) random variable. The notion of dimension here deserves further attention. For one, we have the regressor dimension. This is the dimension of the space of the observed or derived parameters. Not all parameters are identifiable or estimable. More precisely, this dimension is the vector space dimension of the chosen model space. Models are indicated by parameters in this space. The parameter can be unknown or hypothetical. In any case, we have considered the parameters as deterministic. On the other side we have the stochastic component, represented by the error term. In this chapter, we assumed homogeneous errors. So the error term in principle is one dimensional, with distribution taken from a given space of distributions. For the special case of simple Gauss-linear models, the distributions are specified by two parameters, the expected value and the variance. We made the assumption that the model would cover all systematic effects on average, hence the expected value of the error is zero. The variance remains as an unknown nuisance parameter. We avoided the problems coming from this nuisance parameter by restricting ourselves to problems where the nuisance parameter could be replaced by an estimated value and thus be eliminated.

## 2.8 Literature and Additional References

s:s02-bib

[chh92sns] [6] Chambers, J.M.; Hastie, T.J. (eds.) (1992): Statistical Models in S. New York: Chapman & Hall.

[jor93tlm] [15] Jørgensen, B. (1993): The Theory of Linear Models. New York: Chapman & Hall.

[R:Rlang] [21] R Development Core Team (2004): The R language definition.

[gs94nrel1] [26] Sawitzki, G. (1994): Numerical Reliability of Data Analysis Systems. Computational Statistic

tics & Data Analysis 18.2 (1994) 269-286. <<http://www.statlab.uni-heidelberg.de/reports/>>.

- [27] Sawitzki, G. (1994):Report on the Numerical Reliability of Data Analysis Systems. Computational Statistics & Data Analysis/SSN 18.2 (1994) 289-301. <<http://www.statlab.uni-heidelberg.de/reports/>>.

Generated by Sweave from:  
Source: /u/math/j40/cvsroot/lectures/src/SIntro/Rintro/Rnw/S02reg.Rnw.tex,v  
Revision: 1.38  
Date: 2008/06/18 16:48:42  
name:  
Author: j40  
\$Source: /u/math/j40/cvsroot/lectures/src/SIntro/Rintro/Rnw/S02reg.Rnw.tex,v \$  
\$Revision: 1.38 \$  
\$Date: 2008/06/18 16:48:42 \$  
\$name: \$  
\$Author: j40 \$

---

## CHAPTER 3

# Comparison of Distributions

---

*Private note:*

*The chapter title is a misnomer. A better title is needed. This is essentially a second pass through chapter 1.*

*Relations to chapter 1 and chapter 2 needs to be clarified.*

*The discussion data features has been moved here to avoid overloading chapter 1. Its role needs to be clarified.*

ch:03

We begin with the construction of a small gadget that will provide us with example data. The base is a small reaction timer. We draw a “random” point, wait for a mouse click, and record the position of the mouse pointer. To get a stable image for repeated activations, we fix the co-ordinate system.

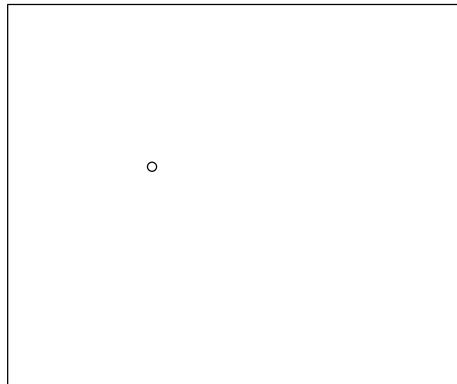
ch03:Sframe03.1

*Example 3.1: Interactive Location*

```
Input  
plot(x = runif(1), y = runif(1),  
      xlim = c(0, 1), ylim = c(0, 1),  
      main = "Please click on the circle",  
      xlab = '', ylab = '',  
      axes = FALSE, frame.plot = TRUE)  
locator(1)
```

```
Output  
$x  
[1] 0.9080969  
  
$y  
[1] 0.3763401
```

Please click on the circle



Now we wrap up the base timer. We record the co-ordinates, try to measure the reaction time, and return the results as a list.

ch03:Sframe03.2

**Example 3.2: Click Timing**

```
click1 <- function(){
  x <- runif(1);y <- runif(1)
  plot(x = x, y = y, xlim = c(0, 1), ylim = c(0, 1),
       main = "Please click on the circle",
       xlab = '', ylab = '',
       axes = FALSE, frame.plot = TRUE)
  clicktime <- system.time(xyclick <- locator(1))
  list(timestamp = Sys.time(),
       x = x, y = y,
       xclick = xyclick$x, yclick = xyclick$y,
       tclick = clicktime[3])
}
```

For later processing we can integrate the list in a `data.frame` and extend this `data.frame` stepwise using `rbind`.

**Example 3.3: Sequential Recording**

```
dx <- as.data.frame(click1())
dx <- rbind(dx, data.frame(click1()))
dx
```

	Output				
	timestamp	x	y	xclick	yclick
elapsed	2008-06-19 09:58:37	0.54538	0.220393	0.90607	0.36551
elapsed1	2008-06-19 09:58:38	0.19637	0.050196	0.90506	0.36551
	tclick				
elapsed	0.455				
elapsed1	0.368				

a:03-click

**Exercise 3.1**

Define a function `click(runs)` that repeats `click1()` a chosen number `runs` of times and returns the result as a `data.frame`. An additional first timing should be considered as a “warming up” and not included in the following evaluation.

Select a number `runs`. Give reasons for your choice of `runs`. Execute `click(runs)` and store the result in a file using `write.table()`. Display the distribution of the component `tclick()` with the methods from chapter [1](#) (distribution function, histogram, box&whisker plot).

### 3.1 Shift/Scale Families, and Stochastic Order

A comparison of distributions can be a demanding task. The mathematical space of distributions is not a number space any more nor a finite dimensional vector space. Properly the room in which distributions reside is a space of probability measures. In simple cases, like for distributions on  $\mathbb{R}$ , we can reduce everything to distribution functions and come to a tractable function space. However, even here a comparison can provide big problems. We do not have any simple order relation to base the comparison upon.

a:03-clickrl

Exercise 3.2	
	<p><a href="#">a:03-click</a></p> <p>Perform exercise 3.1 using the right hand and then again using the left hand. Compare the empirical distributions of the timing data returned by <code>tclick()</code>.</p> <p>The recorded data contain also information about the positions. Define a distance measure <code>dist</code> for the deviation. Give reasons for your definition. Perform a right/left comparison for <code>dist</code>.</p>

We concentrate on the comparison of two distributions only, for example that of the results of two treatments. And we take a simple case: we assume that the observations are independent identically distributed. We use the index notation which is usual for the comparison of treatments in the two sample case.

$Y_{ij}$  independent identically distributed with distribution function  $F_i$

$i = 1, 2$  treatments

$j = 1, \dots, n_i$  observations in treatment group  $i$ .

How do we compare the observations in the treatment groups  $i = 1, 2$ ? The (simple) linear models

$$Y_{ij} = \mu + \alpha_i + \varepsilon_{ij}$$

consider only the case where the difference amounts to a shift  $\Delta = \alpha_1 - \alpha_2$ .

**Notation:** For a distribution with distribution function  $F$  the family

$$F_a(x) = F(x - a)$$

is called the *shift family* for  $F$ .  $a$  is called the shift or location parameter.

Speaking in terms of probabilities, the treatment can shift probability mass in quite different ways than what can be achieved by an additive shift term. We need more general ways to compare distributions. Shift families are not the only framework to consider.

**Notation:** A distribution with distribution function  $F_1$  is *stochastically smaller* than a distribution with distribution function  $F_2$  ( $F_1 \prec F_2$ ), if  $F_1$  takes rather smaller values than as  $F_2$ . This means that  $F_1$  increases sooner.

$$F_1(x) \geq F_2(x) \quad \forall x$$

and

$$F_1(x) > F_2(x) \text{ for at least one } x.$$

For shift families we have: If  $a < 0$ , then  $F_a \prec F$ . The shift results in a parallel shift of the distribution functions.

A typical result of the click experiment (exercise 3.1) is given in figure 3.1. The response times for the right side are stochastically smaller than those for the left side. But the distributions do not belong to a common shift family, because the distribution functions are not parallel.

**ToDo:**

colour for  
left not  
visible in  
b/w

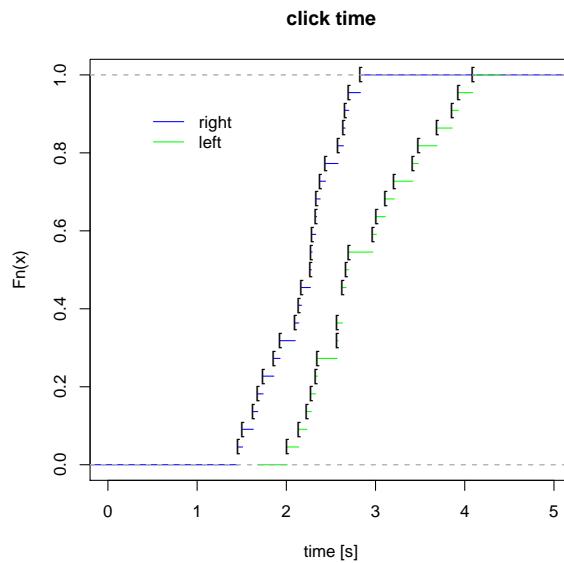


Figure 3.1 Distribution functions for the right/left click time

fig:03clickcdf

a:ch03:1

Exercise 3.3	
	How does a PP plot for $F_1$ against $F_2$ look like if $F_1 \prec F_2$ ? How does a QQ plot for $F_1$ against $F_2$ look like if $F_1 \prec F_2$ ?

Unfortunately the stochastic order defined here is only of a limited value. It does not define a complete order. For shift families it is sufficient. But counter examples can be constructed by minor extensions of shift families.

**Notation:** For a distribution with distribution function  $F$  the family with

$$F_{a,b}(x) = F\left(\frac{x-a}{b}\right)$$

is called the *scale shift family* for  $F$ .

a:ch03:2

Exercise 3.4	
	The scale shift family for the $N(0, 1)$ distribution are the $N(\mu, \sigma^2)$ distributions. Which $N(\mu, \sigma^2)$ distributions are stochastically smaller than the $N(0, 1)$ distribution? Which are stochastically larger? Which distributions have an undefined order relation to $N(0, 1)$ ?

The ordering by scale and location, which comes from the linear theory, and stochastic order are different concepts, and both aspects must be considered separately. Many statistical methods concentrate on aspects which are motivated by scale shift families. Differences going beyond what can be formulated in terms of scale and shift often need special attention.

In chapter 2 we have considered a typical situation for linear models. In principle we dealt with scale shift families. The scale parameter in these models is but a nuisance parameter that can be eliminated. For this, we used an estimator for this stochastic scale parameter, the residual variance, which we cancelled out ultimately. As a particularity for Gauss-linear models we get independent estimators for the expected value and the variance. This helps us in Gauss-linear models to derive statistics which are independent from the nuisance scale parameter.

In general however, we have a spectrum of problems:

- shift alternatives
- shift/scale alternatives
- stochastic ordering
- general alternatives

Testing and estimation methods often concentrate on just one aspect of the problem, the location. The scale parameter of the stochastic error is here but a nuisance parameter. Differences in the shift parameter lead to stochastically monotone relations. Differences in the scale parameter are not categorised so easily. If not only the shift parameter but also the scale parameter may vary, the situation is more complicated and test statistics must first be cleared from the scale as a nuisance parameter.

### 3.2 QQ plot, PP plot, and Comparison of Distributions

s:3.2

As means to compare distributions we already encountered the *PP* plot and the *QQ* plot in chapter 1. In chapter 1 we used these plots to compare an empirical distribution with a theoretical distribution. But of course, with some caution, we can apply it to compare two empirical distributions.

As long as we stay within one shift scale family, the *QQ* plot has at least one advantage over the *PP* plot:

**Remark 3.1** *If  $F_1, F_2$  are distribution functions from a common shift scale family, the *QQ* plot of  $F_1$  against  $F_2$  is a straight line.*

In particular for the Gaussian distributions the *QQ* plot against  $N(0, 1)$  is a valuable tool. Any normal distribution gives a straight line in this plot. The *QQ* plot is already implemented for this situation as function `qqnorm()`.

For the comparison of two samples the corresponding function `qqplot()` can be used: if we denote the empirical quantiles by  $Y_{1,(i:n)}$  resp.  $Y_{2,(i:n)}$ , the plot is the graph  $(Y_{1,(i:n)}, Y_{2,(i:n)})_{i=1\dots n}$  in the case of equal sample sizes. The necessary caution we mentioned above comes in if the sample sizes differ. If the sample sizes differ, R generates nodes by linear interpolation, with the smaller of the two sample sizes determining the number of nodes.

The *PP* plot does not have equivariance properties comparable to those of the *QQ*. If we want to eliminate scale and shift parameters, we have to transform the data first. Mathematical theory however is simple for the *PP* plot. In particular there is a corresponding Kolmogorov-Smirnov test (see section 3.2.1).

On the other hand, the equivariance properties of the *QQ* plots are paid for with structural deficits. In low density regions, few data points control the plot. As a consequence, the plot shows high variance. On the other hand, in these regions quantiles that are near as quantiles are far apart in value space: high variance combines unfavourably with high variability, and correspondingly the *QQ* plot shows large fluctuation. For most text book distributions this means that the *QQ* plot is not very useful in the tails. The *PP* plot does not have this kind of scale deficit, but then it does not have the equivariance properties of the the *QQ* plots. To cope with this, the *PP* plot is generally applied to variable which are standardised in an appropriate way. We will give an example soon.

**help(qqplot)**

<code>qqnorm</code>	<i>Quantile-Quantile Plots</i>
---------------------	--------------------------------

#### Description

`qqnorm` is a generic function the default method of which produces a normal QQ plot of the values in `y`. `qqline` adds a line to a normal quantile-quantile plot which passes through the first and third quartiles.

`qqplot` produces a QQ plot of two datasets.

Graphical parameters may be given as arguments to `qqnorm`, `qqplot` and `qqline`.

#### Usage

```
qqnorm(y, ...)
## Default S3 method:
qqnorm(y, ylim, main = "Normal Q-Q Plot",
      xlab = "Theoretical Quantiles", ylab = "Sample Quantiles",
      plot.it = TRUE, datax = FALSE, ...)

qqline(y, datax = FALSE, ...)

qqplot(x, y, plot.it = TRUE, xlab = deparse(substitute(x)),
       ylab = deparse(substitute(y)), ...)
```

*Arguments*

**x** The first sample for `qqplot`.  
**y** The second or only data sample.  
**xlab, ylab, main** plot labels. The `xlab` and `ylab` refer to the y and x axes respectively if `datax = TRUE`.  
**plot.it** logical. Should the result be plotted?  
**datax** logical. Should data values be on the x-axis?  
**ylim, ...** graphical parameters.

*Value*

For `qqnorm` and `qqplot`, a list with components

**x** The x coordinates of the points that were/would be plotted  
**y** The original y vector, i.e., the corresponding y coordinates *including NAs*.

*References*

Becker, R. A., Chambers, J. M. and Wilks, A. R. (1988) *The New S Language*. Wadsworth & Brooks/Cole.

*See Also*

`ppoints`, used by `qqnorm` to generate approximations to expected order statistics for a normal distribution.

*Examples*

```
require(graphics)

y <- rt(200, df = 5)
qqnorm(y); qqline(y, col = 2)
qqplot(y, rt(300, df = 5))

qqnorm(precip, ylab = "Precipitation [in/yr] for 70 US cities")
```

a:03-clickqq

Exercise 3.5	
	Use the <i>QQ</i> plot to compare the results of the right/left <code>click</code> experiments. Summarise the results.
	(cont.)→

Exercise 3.5	(cont.)
	Combine the right/left <code>tclick</code> data to a vector. Compare the <i>QQ</i> plot with that of Monte Carlo samples taken from the joined vector. <i>Hint:</i> You can draw random samples with <code>sample()</code> . With <code>par(mfrow = c(2, 2))</code> you arrange the display area so that it shows four plots at a time.
**	For <code>sample()</code> use <code>replace = FALSE</code> . How do you have to apply <code>sample()</code> now to split the joint vector into two vectors with Monte Carlo samples? What differences do you expect in comparison to <code>replace = TRUE</code> ?

a:ch03:3

Exercise 3.6	
	Find scale and shift parameters for the right/left <code>click</code> data so that, after using these parameters for transformation, the groups match as good as possible. Describe the differences using these parameters. Use a model formulation in terms of a linear model.
	Use the function <code>boxplot()</code> , to display quartiles and tail behaviour. Compare the information with the information you derived from the scale and shift parameters. <i>Hint:</i> What corresponds to the shift (or location) parameter? What corresponds to the scale parameter?

If representations such as visual representations in displays or numeric representations in summary statistics are affine invariant, scale and shift parameters can be ignored. If representations are not affine invariant, it is often helpful to estimate scale and shift parameters first, then standardise the distributions, and only then to inspect the standardised distributions.

The potential problem with this is that we have to take into account the stochastic behaviour of the scale and shift parameter estimation. The usual way out is to be cautious and use “conservative” tests and robust estimators. The following function tries to transform scale and location to match a standard normal distribution.

---

*Input*

```
ScaleShiftStd <- function (x) {
  xq <- quantile(x[!is.na(x)], c(0.25, 0.75))
  y <- qnorm(c(0.25, 0.75))
  slope <- diff(y)/diff(xq)
  (x-median(x, na.rm = FALSE)) * slope
}
```

---

For direct comparison of distributions, we return to techniques discussed in the first chapter. In principle what has been said there about the comparison to a theoretical distribution can be carried over to the comparison of two distributions, for example to comparison of two treatments. The statistical statements however have to be adapted. In chapter one we had one fixed hypothetically known distribution to be compared with an empirical distribution. This is sometimes called the *one sample case*. Now we are in the *two sample case*. We have to compare

two empirical distributions, with the understanding that the proper targets of comparison are the underlying unknown distributions.

The Monte Carlo band that we used in the one sample case does not have an immediate counterpart. We want to compare two distributions, but we have no distinguished model distribution to draw reference samples from.

However we can modify the idea and construct conditional Monte Carlo bands. Conditioning here means: the construction is dependent on the observed sample values. Assume we have two samples  $Y_{11}, \dots, Y_{1n_1}$  and  $Y_{21}, \dots, Y_{2n_2}$  of independent and, within the group, identically  $F_1$  resp.  $F_2$  distributed observations. If there is no difference between the groups,  $(Y_{11}, \dots, Y_{1n_1}, Y_{21}, \dots, Y_{2n_2})$  is a iid random sample from the common distribution  $F = F_1 = F_2$  with sample size  $n = n_1 + n_2$ . For an iid random sample, each permutation of the indices would arise with equal probability.

This motivates the following procedure: we permute the tuple  $(Y_{11}, \dots, Y_{1n_1}, Y_{21}, \dots, Y_{2n_2})$  and after permutation we assign the first  $n_1$  values to the first group, the remaining values to the second.

As  $n_1, n_2$  increase, the permutation group soon is so big that an exhaustive evaluation becomes impossible. As a way out we use a random sample of permutations. We use the values generated by this procedure to generate Monte Carlo bands.

ToDo:  
ch03: work

out  
ToDo:  
ch03:  
RExercise:  
power

Exercise 3.7	Two Sample Monte Carlo Bands
	Modify the functions for the <i>PP</i> plot and the <i>QQ</i> plot so that Monte Carlo bands for the comparison of two samples are added. (Use a scale/shift standardisation for the <i>PP</i> plot.) <i>Hint:</i> Use the function <code>sample()</code> to generate random permutations.

At larger sample sizes, the mere administrative overhead to generate permutations may become too big. Instead of generating permutations, drawing samples with replacement from the  $n$  values  $(Y_{11}, \dots, Y_{1n_1}, Y_{12}, \dots, Y_{1n_2})$  may be more economical and may be used as a substitute. This approximate solution is known as *bootstrap approximation*<sup>1</sup>.

Since there are only finitely many permutations for a given sample size, for small sample sizes we can enumerate all permutations. We choose bands such that a sufficiently high proportion of all curves are captured within the bands (for example more than 95 %). Permutations which only differ within the groups give the same bands. Hence we need not check all  $n!$  permutations, but only the  $\binom{n}{n_1}$  selections for the assignment to the groups.

a:ch03:6

Exercise 3.8	
**	Augment the <i>PP</i> plot and <i>QQ</i> plot for the <i>click</i> experiments by permutation bands that cover 95 % of the permutations.
	(cont.)→

<sup>1</sup> Caution: there are arbitrarily wild definitions of bootstrap. Always try to get a concise mathematical formulation when speaking about bootstrap.

<b>Exercise 3.8</b>	
*	Generate new plots from the <i>PP</i> -Plots and <i>QQ</i> plots by adding Monte Carlo bands from permutations. Use the envelope of 19 Monte Carlo samples. <i>Hint:</i> Use function <code>sample()</code> to draw a random sample of sample size $n_1$ from $x = (Y_{11}, \dots, Y_{1n_1}, Y_{12}, \dots, Y_{1n_2})$ .
	<i>Hint:</i> See <code>help(sample)</code> .

a:ch03:7

<b>Exercise 3.9</b>	
*	Try to compare the properties of permutation bands, Monte Carlo bands and bootstrap bands on the hypothesis where $F_1 = F_2$ .

If not the distributions, but only single specified parameters are to be compared, an analogous strategy can be used. For example, if we focus on shift alternatives (that is  $F_1$  and  $F_2$  are from a shift family,  $F_1(x) = F_2(x - a)$  for some  $a$ , we can take the mean (or the median) as the parameter of interest. The procedure given above can be used analogously to test the hypothesis that the distributions are not different ( $a = 0$ ), based on the data.

a:ch03:8

<b>Exercise 3.10</b>	
*	Formulate the strategies given above for intervals of single test statistic (example: mean) instead of bands. <i>Hint:</i> Instead of the two mean values for both groups, can you use a single one dimensional statistic?

### 3.2.1 Kolmogorov Smirnov Tests

subs:3.4.2

In chapter 1 we have seen the Kolmogorov-Smirnov test for comparison of a random sample  $(X_i)_{i=1,\dots,n}$  and the corresponding empirical distribution function  $F_n$  with a (fixed, given) distribution  $F$ . The critical test statistic was

$$\sup |F_n - F|.$$

We can modify this test slightly to compare two empirical distributions. Instead of a model distribution  $F$  we now have a second empirical distribution  $G_m$  from observations  $(Y_j)_{j=1,\dots,m}$  based on some underlying (unknown) distribution  $G$ . The critical test statistic is now

$$\sup |F_n - G_m|.$$

The test based on this statistic is in the 2 sample Kolmogorov-Smirnov test which is presented in many text books. This test corresponds to the *PP* plot and allows the construction of bands for the *PP* plot.

We also can use simulation to determine bands. In contrast to the one sample case we do not have a given distribution from which to simulate. Under the hypothesis, that the distributions

$F$  and  $G$  do not differ for independent observations the joined vector  $(X_1, \dots, X_n, Y_1, \dots, Y_m)$  is vector of  $n + m$  independent random numbers with identical distribution  $F = G$ . Given a data set, this relation can be used for simulation. Using a permutation  $\pi$  of the indices from the vector  $Z = (X_1, \dots, X_n, Y_1, \dots, Y_m)$  a new vector  $Z'$  with  $Z'_i = Z_{\pi(i)}$  is generated. The first  $n$  components are used as simulated values  $(X'_i)_{i=1,\dots,n}$ , the remaining  $m$  components as simulated values  $(Y'_j)_{j=1,\dots,m}$ .

a:ch03:14

Exercise 3.11	
*	Implement this algorithm and enhance the $PP$ plot by adding simulated $PP$ plots generated by a for a small number (19?) of permutations.
	Determine the permutation distribution of $\sup  F_n - G_m $ from the simulation and calculate this statistic for the original data. Can you use this comparison to define a test procedure?
	The Kolmogorov-Smirnov test as implemented uses an approximation for the two sample case. In our simulation we know that we simulate under the hypothesis. So any rejection we get is a false rejecton, i.e. an error. Inspect the distribution of the error level under the simulated conditions.

a:ch02:21

Exercise 3.12	
	Use the $QQ$ plot, for a pairwise comparison of the results of the helicopter experiment from chapter 2. Summarise your resultats.

a:ch02:22

Exercise 3.13	
	Inspect the implementation of <code>qqnorm()</code> . Implement an analogous function for the $PP$ plot and apply it to the helicopter data.

### 3.3 Tests for Shift Alternatives

s:3.3

If we make additional distribution assumptions, we can possibly choose better decision rules. For these however the distribution assumptions may be critical. The dependency on these distribution assumptions can be relaxed or avoided, if we can guarantee appropriate assumptions by our experimental procedure. The  $F$ -test which we met in the last chapter is an example of a distribution based method. For the two sample case this test can be re-expressed as a  $t$ -test, which also gives an indication of the direction of the difference. (The square of the  $t$  statistic is an  $F$  statistic.)

---

`help(t.test)`

---

<b>t.test</b>	<i>Student's t-Test</i>
---------------	-------------------------

---

*Description*

Performs one and two sample t-tests on vectors of data.

*Usage*

```
t.test(x, ...)

## Default S3 method:
t.test(x, y = NULL,
       alternative = c("two.sided", "less", "greater"),
       mu = 0, paired = FALSE, var.equal = FALSE,
       conf.level = 0.95, ...)

## S3 method for class 'formula':
t.test(formula, data, subset, na.action, ...)
```

*Arguments*

<b>x</b>	a (non-empty) numeric vector of data values.
<b>y</b>	an optional (non-empty) numeric vector of data values.
<b>alternative</b>	a character string specifying the alternative hypothesis, must be one of "two.sided" (default), "greater" or "less". You can specify just the initial letter.
<b>mu</b>	a number indicating the true value of the mean (or difference in means if you are performing a two sample test).
<b>paired</b>	a logical indicating whether you want a paired t-test.
<b>var.equal</b>	a logical variable indicating whether to treat the two variances as being equal. If TRUE then the pooled variance is used to estimate the variance otherwise the Welch (or Satterthwaite) approximation to the degrees of freedom is used.
<b>conf.level</b>	confidence level of the interval.
<b>formula</b>	a formula of the form <b>lhs ~ rhs</b> where <b>lhs</b> is a numeric variable giving the data values and <b>rhs</b> a factor with two levels giving the corresponding groups.
<b>data</b>	an optional matrix or data frame (or similar: see <b>model.frame</b> ) containing the variables in the formula <b>formula</b> . By default the variables are taken from <b>environment(formula)</b> .
<b>subset</b>	an optional vector specifying a subset of observations to be used.
<b>na.action</b>	a function which indicates what should happen when the data contain NAs. Defaults to <b>getOption("na.action")</b> .
<b>...</b>	further arguments to be passed to or from methods.

*Details*

The formula interface is only applicable for the 2-sample tests.

`alternative = "greater"` is the alternative that `x` has a larger mean than `y`.

If `paired` is `TRUE` then both `x` and `y` must be specified and they must be the same length. Missing values are removed (in pairs if `paired` is `TRUE`). If `var.equal` is `TRUE` then the pooled estimate of the variance is used. By default, if `var.equal` is `FALSE` then the variance is estimated separately for both groups and the Welch modification to the degrees of freedom is used.

If the input data are effectively constant (compared to the larger of the two means) an error is generated.

*Value*

A list with class "`htest`" containing the following components:

<code>statistic</code>	the value of the t-statistic.
<code>parameter</code>	the degrees of freedom for the t-statistic.
<code>p.value</code>	the p-value for the test.
<code>conf.int</code>	a confidence interval for the mean appropriate to the specified alternative hypothesis.
<code>estimate</code>	the estimated mean or difference in means depending on whether it was a one-sample test or a two-sample test.
<code>null.value</code>	the specified hypothesized value of the mean or mean difference depending on whether it was a one-sample test or a two-sample test.
<code>alternative</code>	a character string describing the alternative hypothesis.
<code>method</code>	a character string indicating what type of t-test was performed.
<code>data.name</code>	a character string giving the name(s) of the data.

*See Also*

`prop.test`

*Examples*

```
require(graphics)

t.test(1:10,y=c(7:20))      # P = .00001855
t.test(1:10,y=c(7:20, 200)) # P = .1245    -- NOT significant anymore

## Classical example: Student's sleep data
plot(extra ~ group, data = sleep)
## Traditional interface
with(sleep, t.test(extra[group == 1], extra[group == 2]))
## Formula interface
t.test(extra ~ group, data = sleep)
```

The text book derivation for the  $t$ -test supposes, that we have independent identically distributed samples from normal distributions. In fact, weaker assumptions are sufficient. Writing the  $t$ -test statistic as

$$t = \frac{\widehat{\mu}_1 - \widehat{\mu}_2}{\sqrt{Var(\widehat{\mu}_1 - \widehat{\mu}_2)}} \quad (3.1)$$

eq:03tstat

we see that  $t$  is  $t$ -distributed if  $\widehat{\mu}_1 - \widehat{\mu}_2$  has a normal distribution,  $(Var(\widehat{\mu}_1 - \widehat{\mu}_2))$  has a  $\chi^2$  distribution, and both terms are independent. The central limit theorem guarantees that under mild conditions  $\widehat{\mu}_1 - \widehat{\mu}_2$  at least asymptotically has a normal distribution. Similarly under not too strict conditions  $(Var(\widehat{\mu}_1 - \widehat{\mu}_2))$  behaves asymptotically as  $\chi^2$  distributed. If both terms are (approximatively) independent,  $t$  is approximatively  $t$  distributed.

a03-tsimg

Exercise 3.14	
*	<p>Use a simulation to inspect the distribution of <math>\bar{Y}</math>, <math>\widehat{Var(Y)}</math> and the <math>t</math> statistic for <math>Y</math> from a uniform distribution <math>U[0, 1]</math> with sample size <math>n = 1, \dots, 10</math>. Compare the distributions from the simulation with the corresponding normal-, <math>\chi^2</math>- resp. <math>t</math> distribution.</p> <p>Use a simulation to inspect the distribution of <math>\bar{Y}</math>, <math>\widehat{Var(Y)}</math> and the <math>t</math> statistic for <math>Y</math> from a mixture, consisting at 90% from a <math>N(0, 1)</math>- and at 10% from a <math>N(0, 10)</math> distribution, with sample size <math>n = 1, \dots, 10</math>. Compare the distributions from the simulation with the corresponding normal-, <math>\chi^2</math>- resp. <math>t</math> distribution.</p>

The  $t$ -test has some robustness which can give it approximate validity outside its normal model context. There are however ways to free us completely from the normal distribution assumption. If we proceed as in the  $F$ -test resp.  $t$ -test, but use ranks instead of the original data, we get test procedures that are distribution independent (at least as long as no ties occur). The Wilcoxon test is a distribution independent variant of the  $t$ -test. Theoretically it corresponds exactly to a  $t$ -test, applied to the (jointly) rank transformed data. Like the  $t$ -test, this test is only designed to test the null hypothesis (no difference) against a shift alternative. For practical application, arithmetic simplifications can be used. This may hide the relation between the usual formula for the  $t$ -test and for the Wilcoxon test.

To apply the Wilcoxon test, on the one hand the test statistic has to be calculated. On the other hand, to determine the critical values, the distribution function has to be evaluated. If all observations are distinct, this function depends only on  $n_1$  and  $n_2$ , and fairly simple algorithms are available. These are provided in the R base and used by `wilcox.test()`. However if there are ties in the data, that is there are values occurring more than once, the distribution depends on the special pattern of these ties and the calculation is laborious. `wilcox.test()` returns to approximations in this case. For an exact evaluation (in contrast to approximative), the necessary algorithms are available as well. To use them, you need `library(coin)`. The exact variant of the Wilcoxon tests for example is implemented as `wilcox_test()`.

It is a classical branch in statistics to characterise distribution independent methods based on ranks and to compare them with independent methods based on Monte Carlo calculations and their variations. You find literature for this area under the key words “rank tests” and “distribution free methods”. Additional R functions are provided in `library(coin)` as well as in some specialised libraries.

**ToDo:**  
add power  
function

**ToDo:**  
clarify: in-  
formation  
 $<- >$  as-  
sumptions

**ToDo:**  
RB: If you  
believe  
the dis-  
tribution  
assump-  
tion, for  
some good  
reason

Of course information loss is a concern. If we restrict to the data and make no or weak assumptions on the distributions, we have less information than in a model with explicit distribution assumptions. If we reduce the data to ranks, we may give away additional information. This loss in information can be measured for instance in terms of asymptotic relative efficiency, that is the (asymptotic) sample size proportion which is needed to get a comparative power using the competing test. Under the assumption of normality, the asymptotic relative efficiency of Wilcoxon test is 94%. So if the assumption of a normal distribution holds, the (optimal)  $t$  only needs 94% of the sample size needed by the Wilcoxon test. 6% of the sample size are the costs for the reduction to ranks. But if the normal distribution does not apply, the  $t$ -test possibly can break down, while the Wilcoxon test stays valid for the shift alternative.

**ToDo:**

add power  
simulation - see  
Rnw.tex  
source

---

**help(wilcox.test)**

---

**wilcox.test**

*Wilcoxon Rank Sum and Signed Rank Tests*

---

*Description*

Performs one and two sample Wilcoxon tests on vectors of data; the latter is also known as ‘Mann-Whitney’ test.

*Usage*

```
wilcox.test(x, ...)

## Default S3 method:
wilcox.test(x, y = NULL,
            alternative = c("two.sided", "less", "greater"),
            mu = 0, paired = FALSE, exact = NULL, correct = TRUE,
            conf.int = FALSE, conf.level = 0.95, ...)

## S3 method for class 'formula':
wilcox.test(formula, data, subset, na.action, ...)
```

*Arguments*

<b>x</b>	numeric vector of data values. Non-finite (e.g. infinite or missing) values will be omitted.
<b>y</b>	an optional numeric vector of data values.
<b>alternative</b>	a character string specifying the alternative hypothesis, must be one of “ <b>two.sided</b> ” (default), “ <b>greater</b> ” or “ <b>less</b> ”. You can specify just the initial letter.
<b>mu</b>	a number specifying an optional parameter used to form the null hypothesis. See ‘Details’.
<b>paired</b>	a logical indicating whether you want a paired test.
<b>exact</b>	a logical indicating whether an exact p-value should be computed.

<code>correct</code>	a logical indicating whether to apply continuity correction in the normal approximation for the p-value.
<code>conf.int</code>	a logical indicating whether a confidence interval should be computed.
<code>conf.level</code>	confidence level of the interval.
<code>formula</code>	a formula of the form <code>lhs ~ rhs</code> where <code>lhs</code> is a numeric variable giving the data values and <code>rhs</code> a factor with two levels giving the corresponding groups.
<code>data</code>	an optional matrix or data frame (or similar: see <code>model.frame</code> ) containing the variables in the formula <code>formula</code> . By default the variables are taken from <code>environment(formula)</code> .
<code>subset</code>	an optional vector specifying a subset of observations to be used.
<code>na.action</code>	a function which indicates what should happen when the data contain NAs. Defaults to <code>getOption("na.action")</code> .
<code>...</code>	further arguments to be passed to or from methods.

### Details

The formula interface is only applicable for the 2-sample tests.

If only `x` is given, or if both `x` and `y` are given and `paired` is `TRUE`, a Wilcoxon signed rank test of the null that the distribution of `x` (in the one sample case) or of `x - y` (in the paired two sample case) is symmetric about `mu` is performed.

Otherwise, if both `x` and `y` are given and `paired` is `FALSE`, a Wilcoxon rank sum test (equivalent to the Mann-Whitney test: see the Note) is carried out. In this case, the null hypothesis is that the distributions of `x` and `y` differ by a location shift of `mu` and the alternative is that they differ by some other location shift (and the one-sided alternative "`greater`" is that `x` is shifted to the right of `y`).

By default (if `exact` is not specified), an exact p-value is computed if the samples contain less than 50 finite values and there are no ties. Otherwise, a normal approximation is used.

Optionally (if argument `conf.int` is true), a nonparametric confidence interval and an estimator for the pseudomedian (one-sample case) or for the difference of the location parameters `x-y` is computed. (The pseudomedian of a distribution  $F$  is the median of the distribution of  $(u+v)/2$ , where  $u$  and  $v$  are independent, each with distribution  $F$ . If  $F$  is symmetric, then the pseudomedian and median coincide. See Hollander & Wolfe (1973), page 34.) If exact p-values are available, an exact confidence interval is obtained by the algorithm described in Bauer (1972), and the Hodges-Lehmann estimator is employed. Otherwise, the returned confidence interval and point estimate are based on normal approximations.

With small samples it may not be possible to achieve very high confidence interval coverages. If this happens a warning will be given and an interval with lower coverage will be substituted.

### Value

A list with class "`htest`" containing the following components:

<code>statistic</code>	the value of the test statistic with a name describing it.
<code>parameter</code>	the parameter(s) for the exact distribution of the test statistic.
<code>p.value</code>	the p-value for the test.

<code>null.value</code>	the location parameter <code>mu</code> .
<code>alternative</code>	a character string describing the alternative hypothesis.
<code>method</code>	the type of test applied.
<code>data.name</code>	a character string giving the names of the data.
<code>conf.int</code>	a confidence interval for the location parameter. (Only present if argument <code>conf.int = TRUE</code> .)
<code>estimate</code>	an estimate of the location parameter. (Only present if argument <code>conf.int = TRUE</code> .)

### Warning

This function can use large amounts of memory and stack (and even crash R if the stack limit is exceeded) if `exact = TRUE` and one sample is large (several thousands or more).

### Note

The literature is not unanimous about the definitions of the Wilcoxon rank sum and Mann-Whitney tests. The two most common definitions correspond to the sum of the ranks of the first sample with the minimum value subtracted or not: R subtracts and S-PLUS does not, giving a value which is larger by  $m(m + 1)/2$  for a first sample of size  $m$ . (It seems Wilcoxon's original paper used the unadjusted sum of the ranks but subsequent tables subtracted the minimum.)

R's value can also be computed as the number of all pairs  $(x[i], y[j])$  for which  $y[j]$  is not greater than  $x[i]$ , the most common definition of the Mann-Whitney test.

### References

- David F. Bauer (1972), Constructing confidence sets using rank statistics. *Journal of the American Statistical Association* **67**, 687–690.
- Myles Hollander & Douglas A. Wolfe (1973), *Nonparametric Statistical Methods*. New York: John Wiley & Sons. Pages 27–33 (one-sample), 68–75 (two-sample).
- Or second edition (1999).

### See Also

- `psignrank`, `pwilcox`.
- `wilcox.exact` in `exactRankTests` covers much of the same ground, but also produces exact p-values in the presence of ties.
- `wilcox_test` in package `coin` for exact and approximate *conditional* p-values for the Wilcoxon tests.
- `kruskal.test` for testing homogeneity in location parameters in the case of two or more samples; `t.test` for an alternative under normality assumptions [or large samples]

*Examples*

```

require(graphics)
## One-sample test.
## Hollander & Wolfe (1973), 29f.
## Hamilton depression scale factor measurements in 9 patients with
## mixed anxiety and depression, taken at the first (x) and second
## (y) visit after initiation of a therapy (administration of a
## tranquilizer).
x <- c(1.83, 0.50, 1.62, 2.48, 1.68, 1.88, 1.55, 3.06, 1.30)
y <- c(0.878, 0.647, 0.598, 2.05, 1.06, 1.29, 1.06, 3.14, 1.29)
wilcox.test(x, y, paired = TRUE, alternative = "greater")
wilcox.test(y - x, alternative = "less")      # The same.
wilcox.test(y - x, alternative = "less",
            exact = FALSE, correct = FALSE) # H&W large sample
                                              # approximation

## Two-sample test.
## Hollander & Wolfe (1973), 69f.
## Permeability constants of the human chorioamnion (a placental
## membrane) at term (x) and between 12 to 26 weeks gestational
## age (y). The alternative of interest is greater permeability
## of the human chorioamnion for the term pregnancy.
x <- c(0.80, 0.83, 1.89, 1.04, 1.45, 1.38, 1.91, 1.64, 0.73, 1.46)
y <- c(1.15, 0.88, 0.90, 0.74, 1.21)
wilcox.test(x, y, alternative = "g")          # greater
wilcox.test(x, y, alternative = "greater",
            exact = FALSE, correct = FALSE) # H&W large sample
                                              # approximation

wilcox.test(rnorm(10), rnorm(10, 2), conf.int = TRUE)

## Formula interface.
boxplot(Ozone ~ Month, data = airquality)
wilcox.test(Ozone ~ Month, data = airquality,
            subset = Month %in% c(5, 8))

```

a:ch03:11

Exercise 3.15	
	Use the Wilcoxon test, to compare the results of the right/left <i>click</i> experiment.

a:ch03:12

Exercise 3.16	
***	<p>In the the right/left <i>click</i> experiment several effects contribute to the response time. Some problems:</p> <ul style="list-style-type: none"> <li>• The response time comprises reaction time, time for the large scale movement of the mouse, time for fine adjustment etc.</li> <li>• For the right/left movement in general a swivel of the hand is sufficient. For forward/backward movement in general a movement of the arm is necessary. It is not to be expected that both movements have a comparable statistical behaviour.</li> <li>• Subsequent record my be affected by a training effect, or by a tiring effect.</li> </ul> <p>Can you modify the experiment or the evaluation so that differences in the reaction time component can be investigated?</p> <p>Can you modify the experiment or the evaluation so that differences in the precision of the position of the click can be investigated?</p>
***	<p>Inspect and document for yourself the right/left differences in reaction time and precision. Summarise your results as a report..</p>

a:ch03:13

Exercise 3.17	
	<p>Use the shift/scale families of <math>N(0, 1)</math> and <math>t(3)</math> and design a setting to compare the performance of the Wilcoxon test with that of the <math>t</math>-test for each of these families.</p> <p>Perform the comparison in a simulation with sample sizes <math>n_1 = n_2 = 10, 20, 50, 100</math> and summarise your results.</p> <p>Do an analogous comparison using simulation data from the lognormal distribution.</p>

### 3.4 A Road Map

If you go hiking, a road map is on the middle of the way, and the most important information is a point saying “You Are Here”. We are in a similar situation. At this point, we can outline the road we have taken, and the perspectives ahead.

We started with basic data analysis in chapter [I](#). This chapter discussed independent, identically distributed data. However the discussion was still in a void space, missing a clear target.

The linear models in chapter [II](#) gave a well defined framework, and put chapter [I](#) in its place. We need basic data analysis, but if we follow the approach of linear models we need diagnostics based on residuals. For linear models, the stochastic part is in the error term. But the stochastic errors are not directly observable. We have to resume on indirect inference, based on the residuals. We have two ways we can go. We can derive diagnostic indices based on the residuals, which are at least approximatively independent, identically distributed data and use the basic analysis as in chapter [I](#). Or we can modify the basic analysis to cope with the dependent inhomogeneous structure derivable from the data. Both approaches are feasible. And both approaches have to

cope with the problem that the residuals are affected both by the data and by our estimation step. So we have to avoid running into circles, such as hidden outliers or masking effects.

In this chapter, we have put linear models in their place. For the simple two sample situation, linear models boil down to comparison of two treatments. What they provide is analysis for a shift. Formally, linear models consider a shift/scale family, but at least for simple linear models the scale is just a nuisance parameter which is removed as soon as possible. If you have some good reasons to believe in the distribution assumptions, linear models are a first choice. More often, you may be working in a field where experimental setups and appropriate scale transformations are well known enough to at least believe in the distribution assumptions with good confidence. (Be aware of the pitfalls. The assumptions, or transformations to enforce the assumptions may be well established for a “normal” population. But are they established for a population to compare with, for example under treatment?) If the assumptions can be made: fine, but use diagnostics as a second check). If they are in doubt, you can use non-parametric methods at a negligible cost, keeping competitive performance when the distribution assumptions are satisfied, but guaranteeing quality even if the distribution assumptions fail. However the standard methods are only tuned for shift alternatives.

The main message of this chapter is that the world is richer than what shift alternatives suggest. Be prepared for the unexpected.

The other open horizon is to go to higher dimensions. So far, the data were one dimensional. Well, not quite. The data in chapter 2 were essentially  $p + 1$  dimensional. It was the specific view of the regression model which factorized the data into a regressor part  $X$  of dimension  $p$  which considered non-stochastic, and a response part  $Y$  which was modeled as stochastic, and of course it is the question to the application whether this split is adequate, or whether the data should be analysed as a joined higher dimensional data set. Dimensionality is a topic still to come. This topic will be touched in the next chapter.

With this said, we can return to our well known procedures. We will have a closer look at their power, and see how this information is related to practical planning of experiments.

### 3.5 Power and Confidence

#### 3.5.1 Theoretical Power and Confidence

We can use the  $t$ -test to illustrate how tests are constructed in general. The test uses a test statistic, here the  $t$ -test statistic for comparison of two groups (3.1). We know the distribution of this statistic: in case of the  $t$ -test, given independent errors with normal distribution and equal variance the test statistic has a  $t(n_1 + n_2 - 2)$  distribution. For given level  $\alpha$  we can use the distribution function to read off limits which are exceeded resp. missed only with a probability of  $\alpha$ . If we use both limits we get a two sided region with error probability  $2\alpha$ .

**To Do:**  
Simulation,  
Power

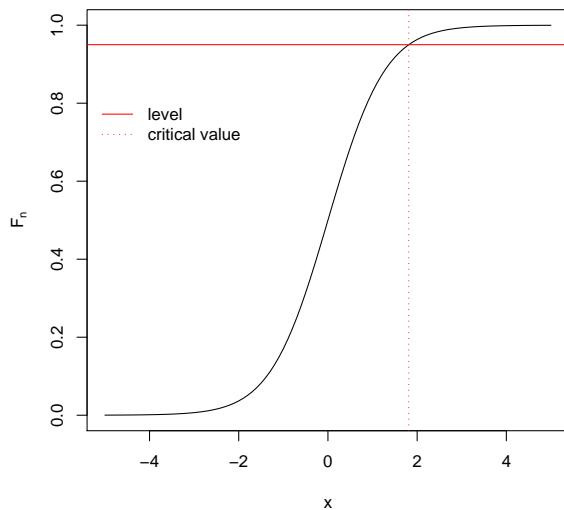
fig:ch03testcrit

**Example 3.4: Test Construction***Input*

```

n1<- 6; n2 <- 6
df <- n1 + n2 -2
alpha <- 0.05
curve(pt(x,df=df),from=-5, to=5, ylab= expression(F[n]),
      main="t-test: critical value")
abline(h=1-alpha, col="red")      # cut at upper quantile
abline(v=qt(1-alpha, df=df), lty=3, col="red") # get critical value
legend("topleft", legend=c("level","critical value"),
      lty=c(1,3),col="red",
      bty="n", inset=c(0,0.2))

```

**t-test: critical value**

If, for example, we want to test the hypothesis  $\mu_1 = \mu_2$  against the alternative  $\mu_1 > \mu_2$ , we use the region above the upper of these limits as rejection region. We know that if the hypothesis applies, a random test statistic will fall into this region with probability of at most  $\alpha$ .

**To Do:**  
introduce  
noncentral  
t in ch. 2

For the  $t$ -test we know more. Under the model assumptions of independent errors with normal distribution and equal variance, the  $t$ -test statistic always has a  $t$  distribution. On the hypothesis it has a  $t$  distribution with non centrality parameter 0, that is it has a central  $t$  distribution. On the alternative where  $\mu_1 \neq \mu_2$ , but still for normal distribution with equal variance, we have a  $t$  distribution with non centrality parameter  $(\mu_1 - \mu_2)\sigma^{-1}\sqrt{n_1 n_2 / (n_1 + n_2)}$ . Under model assumptions this gives the power of the test for each alternative, that is, the probability to reject the hypothesis if a specific alternative applies.

**To Do:**  
ig:ch03test what  
one sided  
or two  
sided?

**Example 3.5: Power**

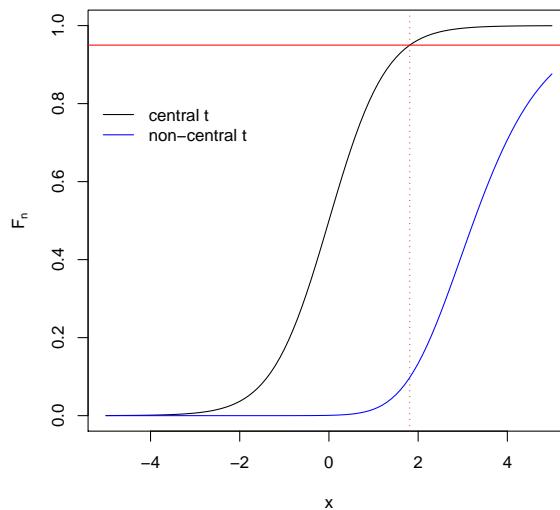
*Input*

```

n1<- 6; n2 <- 6
df <- n1 + n2 -2
alpha <- 0.05
curve(pt(x,df=df),from=-5, to=5, ylab= expression(F[n]),
      main="central and non central t-distribution")
abline(h=1-alpha, col="red")      # cut at upper quantile
abline(v=qt(1-alpha, df=df), lty=3, col="red") # get critical value
n1 <- 5
n2 <- 5
n <- n1+n2
theta <- 2
ncp <- theta * sqrt(n1 * n2/(n1+n2))
curve(pt(x,df=df, ncp=ncp),add=TRUE, col="blue")
legend("topleft", legend=c("central t","non-central t"),
       lty=1, col=c("black","blue"),
       bty="n", inset=c(0,0.2))

```

central and non central t-distribution



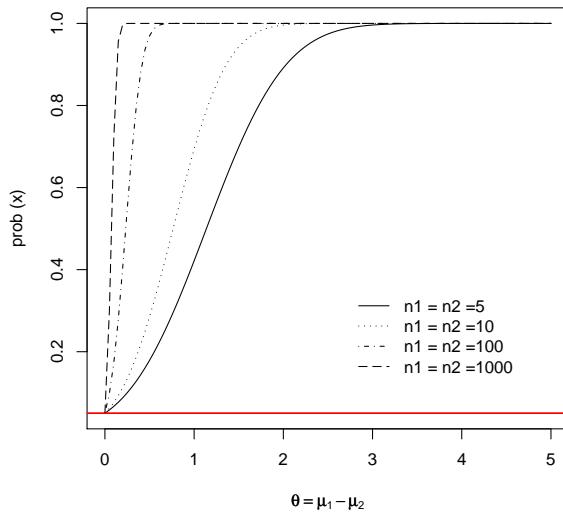
We can display the power of the test by plotting the rejection probability as a function of  $(\mu_1 - \mu_2)\sigma^{-1}$ .<sup>2</sup>

<sup>2</sup> Conventionally the term “power function” is only used for the rejection probability on the alternative, that is for example for  $(\mu_1 - \mu_2) > 0$ .

*Example 3.6: Power Function*

```
Input
tpower <- function(n1, n2, alpha,...){
  df <- n1 + n2 -2
  tlim <- qt(1-alpha,df=df)
  prob <- function(theta){
    pt(tlim, df = df,
       ncp = theta * sqrt(n1 * n2/(n1+n2)),
       lower.tail=FALSE)
  }
  curve(prob, 0, 5, xlab=expression(theta==mu[1]-mu[2]), ...)
  abline(h=alpha, col="red")
}
tpower(5, 5, 0.05, main="power function for selected sample sizes")
tpower(10, 10, 0.05, add =TRUE, lty = 3)
tpower(100,100, 0.05, add =TRUE, lty = 4)
tpower(1000, 1000, 0.05, add =TRUE, lty = 5)
legend("bottomright",
       lty=c(1,3,4,5),
       legend=c("n1 = n2 =5", "n1 = n2 =10",
               "n1 = n2 =100", "n1 = n2 =1000"),
       inset=0.1, bty="n")
```

power function for selected sample sizes



This relation can also be used to calculate the sample size needed to reject the hypothesis with a probability  $\alpha$  if it applies (false rejection), but to reject the hypothesis if some specified alternative applies (power or correct rejection). Function `power.t.test()` is already prepared to do this calculation.

ch03powert

*Example 3.7: Sample Size Determination*

*power.t.test(delta=2,*

*power=0.8,*

*sig.level=0.01,*

*type="two.sample",*

*alternative="one.sided")*

*Input*

*Two-sample t test power calculation*

*n = 6.553292*

*delta = 2*

*sd = 1*

*sig.level = 0.01*

*power = 0.8*

*alternative = one.sided*

*Output*

*NOTE: n is number in \*each\* group*

*3.5.2 Simulated Power and Confidence*

If the theoretical properties of a test are known, this is the best way to analyse its properties. In an environment as R we have the possibility to assess the power even if theoretical results are not available or not accessible. For chosen alternatives we can generate random samples, perform tests and find the relative proportion of samples that lead to a rejection. If we generate *nsimul* independent random samples with identical distribution, the number of rejections is a random variable with binomial distribution and

$$\hat{p} = \frac{\#rejections}{nsimul}$$

is an estimator for the rejection probability.

As an example we see how the *t*-test performs if the data have a log-normal distribution. We compare two groups, each with sample size  $n_1 = n_2 = 10$ . We generate data under the hypothesis first:

ch03tlognH

*Example 3.8: Violation of Assumptions (log-normal)*


---

*Input*

```
nsimul <- 300
n1<- 10; n2 <- 10
alpha <- 0.01
x <- 0
for (i in 1:nsimul) {
  if (t.test(exp(rnorm(n1)),exp(rnorm(n2)),
  alternative="less",
  var.equal = TRUE)$p.value < alpha){
    x <- x+1}
}
p <- x/nsimul
cat("estim p", p)
```

---

*Output*

```
estim p 0.006666667
```

The function `prop.test()` does not only calculate the estimator for the binomial fraction, but also a confidence set.

*Example 3.9: Binomial Confidence Set*


---

*Input*

```
prop.test(n=nsimul, x=x)
```

---

*Output*

```
1-sample proportions test with continuity correction

data: x out of nsimul, null probability 0.5
X-squared = 290.0833, df = 1, p-value < 2.2e-16
alternative hypothesis: true p is not equal to 0.5
95 percent confidence interval:
 0.001155561 0.026512753
sample estimates:
 p
0.006666667
```

Analogous for example under the alternative with  $\log(x_2)$  distributed as  $N(\mu_2, 1)$  with  $\mu_2 = 1$ .

**Example 3.10:**

<pre> nsimul &lt;- 300 n1&lt;- 10; n2 &lt;-10 alpha &lt;- 0.01 x&lt;-0 for (i in 1:nsimul) {   if (t.test(exp(rnorm(n1)),exp(rnorm(n2, mean = 1)),     alternative="less",     var.equal = TRUE)\$p.value &lt; alpha){     x &lt;- x+1} } p &lt;- x/nsimul cat("estim p", p) </pre>	<i>Input</i>
---	--------------

<pre>estim p 0.1366667</pre>	<i>Output</i>
------------------------------	---------------

<pre> prop.test(n = nsimul, x = x) </pre>	<i>Input</i>
---	--------------

<pre> 1-sample proportions test with continuity correction  data: x out of nsimul, null probability 0.5 X-squared = 156.9633, df = 1, p-value &lt; 2.2e-16 alternative hypothesis: true p is not equal to 0.5 95 percent confidence interval:  0.1009148 0.1819944 sample estimates:       p 0.1366667 </pre>	<i>Output</i>
---	---------------

In `library(binom)` several tools are provided for a differentiated analysis of the binomial distribution.

The confidence intervals in this example show that a simulation sample size of  $nsimul = 300$  only gives rough results. For simulations, we want better control over the simulation precision. We can choose the simulation sample size to improve the precision as we like. A sample size calculation can be done with `power.prop.test()`. In experiments, an exact sample size calculation may be necessary. For simulations, often an estimation of the sample size is sufficient.

With  $\hat{p} := Z/n$  as estimator for a probability  $p$  we have  $E(\hat{p}) = p$  and  $Var(\hat{p}) = p(1-p)/n$ . If  $p$  is actually the value of the parameter, an error has to be judged relative to this target parameter. At  $p = 50\%$  an error of  $\pm 1\%$  has to be judged differently than at  $p = 99\%$ . The relative error, the *coefficient of variation*, is

$$\frac{\sqrt{Var(\hat{p})}}{E(\hat{p})} = \sqrt{\frac{1-p}{np}}.$$

To get a coefficient of variation of at most  $\eta$ , we need a sample size

$$n \geq \frac{1-p}{p\eta^2}.$$

If  $n$  and  $p$  are in an order of magnitude where a normal approximation applies, we have an approximate confidence interval with limits

$$\hat{p} \pm \phi_{1-\alpha/2} \sqrt{\frac{\hat{p}(1-\hat{p})}{n}} \text{ for a confidence level } 1-\alpha.$$

If the length of the confidence interval should not exceed  $\eta p$ , we need a sample size

$$n \geq \frac{\phi_{1-\alpha/2}^2(1-p)}{p\eta^2}.$$

As usual, the choice of  $\alpha$  is up to us. For example for  $\alpha = 1\%$  with  $\phi_{1-\alpha/2} = 2.575829$  we get the values in table 3.28. If we work with higher quantiles, we will try to bound the error relative to  $1-p$ . Examples are in table 3.28.

Table 3.28 Required sample size for two-sided confidence intervals with relative length  $\leq \eta$

$p$	$1-p$	$n(\alpha = 10\%)$		$n(\alpha = 1\%)$	
		$\eta = 0.1$	$\eta = 0.01$	$\eta = 0.1$	$\eta = 0.01$
0.500	0.500	271	27 055	663	66 349
0.250	0.750	812	81 166	1990	199 047
0.100	0.900	2435	243 499	5971	597 141
0.010	0.990	26 785	2 678 488	65 685	6 568 547
0.001	0.999	270 283	27 028 379	662 826	662 822 617

tab:simulneta

It is worth to memorise the rough numbers: To estimate a probability in the range of magnitude of  $50\% \pm 5\%$  with 90% confidence, about 300 simulations are necessary. To estimate a value of about 99% up to  $\pm 0.1\%$ , we need 30000 simulations.

### 3.5.3 Quantile Estimation by Simulation

The other side of the problem above is to estimate a quantile from a random sample. We already know that, for a random variable  $X$  with continuous distribution function  $F$ , the variable  $F(X)$  has a uniform distribution on  $[0, 1]$ . For the quantile estimation we need the distribution function, evaluated at the order statistics. We already met this in chapter ??, where we had theorem ??.

thm:simulF

**Theorem 3.2** For independent observations  $X_i, i = 1, \dots, n$  from a continuous distribution function  $F$ , let  $X_{(k:n)}$  denote the  $k$ -th order statistic. Then

$$F(X_{(k:n)}) \sim P_{beta}(\cdot; k, n-k+1).$$

We repeat:

**Remark 3.3** In general the beta distribution is skewed. The expected value of the  $\text{beta}(k, n - k + 1)$  distribution is  $k/(n + 1)$ . For an unbiased estimation of the quantile  $x_p$  use  $X_{(k:n)}$  with  $k/(n + 1) = p$ . The “plug in” approximation  $k/n = p$  gives a biased estimation.

The theorem can be applied directly to get an upper or lower estimate for quantiles. As a special case we can try to use the minimum of the observed values  $X_{(1:n)}$  as lower estimate for the  $p$  quantile. The confidence level is

$$P(X_{(1)} \leq F_p) = P(F(X_{(k)}) \leq p) = I_p(1, n),$$

where  $I$  is the incomplete beta integral. For the special parameters  $(1, n)$  the beta density simplifies to  $n(1 - p)^{n-1}$  and we get the incomplete beta integral  $I_p(1, n) = 1 - (1 - p)^n$ . Hence

$$P(X_{(1)} \leq F_p) = 1 - (1 - p)^n$$

and we can guarantee a confidence level  $1 - \alpha$ , if

$$n \geq \frac{\ln \alpha}{\ln(1 - p)}.$$

Similarly, the observed maximum value can be used as an upper estimate for the  $p$ -quantile. By symmetry we get a confidence level of  $1 - \alpha$ , if

$$n \geq \frac{\ln \alpha}{\ln p}.$$

Examples are given in table 3.29. tab:simulnquant

Table 3.29 Required sample size for estimation of a quantile with confidence level  $\geq 1 - \alpha$

p		n			
$X_{(1)} \leq F_p$	$X_{(n)} \geq F_p$	$\alpha = 10\%$	$\alpha = 5\%$	$\alpha = 1\%$	$\alpha = 0.5\%$
0.500	0.500	4	5	7	8
0.250	0.750	9	11	17	19
0.100	0.900	22	29	44	51
0.010	0.990	230	299	459	528
0.001	0.999	2302	2995	4603	5296

Again it is worth to memorise the rough numbers: to get a one sided estimate for a 1% (99%)-quantile of a continuous distribution function with a confidence of 99%, about 500 simulations are needed.

We can combine one sided limits to intervals. The corresponding result to calculate the probability of intervals is in corollary ???: cor:beta2s

**Corollary 3.4** For the  $k_1$ th and  $k_1 + k_2$ th order statistic, the interval  $(X_{(k_1:n)}, X_{(k_1+k_2:n)})$  is a confidence interval for the  $p$  quantile with coverage probability

$$I_p(k_1, n - k_1 + 1) - I_p(k_1 + k_2, n - k_1 - k_2 + 1).$$

cor:sumul2s

The simulation sample sizes for quantile estimation are considerably smaller than those needed for the estimation of comparable probabilities. In retrospect, this is not surprising: the question whether an observation exceeds some specific quantile is simpler than the task to estimate the  $p$ -value. The sample size can often reduce even more drastically if the question at hand is reduced to a test problem.

Without additional distribution assumptions this gives a first possibility to choose the size of a simulation. In special situations, cute ideas can allow for a drastic reduction of the simulation sample size. But for a start, the calculations given above provide a basis for determining the simulation sample size.

### 3.6 Qualitative Features of Distributions

s:03features

At this point, we return to ask for the purpose of data analysis. Chapter <sup>ch:01</sup> was quite open to question on this. We met some traditional aspect, such as getting information about quantiles or moments. From a more general point of view, we met questions about location and scale. Moments or quantiles found their place as intermediates to derive this information of interest.

<sup>ch:01</sup> Chapter <sup>ch:01</sup> lead to a more focussed view. On the one hand, specific parameters such as the regression coefficients and the scale parameter of the error distribution appeared as central target. On the other hand, indicators for deviations from the linear model or trouble makers such as effective outlier were of interest. This is a typical application example which leads to an analysis of the data which must be linked to the intended statistical analysis of the data, and hence to the intended evaluation method.

Focus has changed with available technology. In the first half of the last century, linear model were the prominent method which was could be covered computationally, and the needed recipe was simple. First, try to get hold of the first two moments (mean and variance of your data). Second, try to find a transformation such that in the framework of a linear model the next two moments of the estimated error are approximatively as those of a normal distribution. So the third moment, the skewness, should be zero after transformation. The fourth was compared to that of a normal distribution, yielding a kurtosis excess, which should be transformed to zero. Then a linear model would be applied.

Technology has advanced, more computing power is available, and occassionally we have more advanced algorithms. So this “method of the moments” has implicitly migrated into estimations and lost importance.

We have a remainder of this point of view in the box&whisker plot. But it goes beyond the traditional view. It uses the median as a location estimator, but has the upper and lower quartile yielding possibly unsymmetric scale estimators to the upper and lower range. The information about out and far out points is an additional information which was not considered in earlier approaches.

On the other side, we have what can be considered density estimators. The histograms are early representants of this family. Kernel density estimators are followers. Both seem to give an intuitive impression of the distribution. But this view changes with different settings. There is a bandwidth, or smoothing parameter, which critically effects the impression. Both do not easily yield for a multi scale representation.

The problem does not seem to be on the side of the method, but on the application side. Which features are relevant for the application?

Theory can deliver solutions for standard problems, such as characterisation of location and scale. This may be helpful to have a comparison basis. Now assume that location and scale are analysed and have been adjusted. What is the next step for comparison?

If we have a shift/scale family, all is solved now. If we have a stochastic ordering, comparison of distribution functions may help. But a typical example may just shift a part of the distribution.

A particular case may arrive where you have a distinct group of responders. So in the two sample case, a possibly unimodal distribution in the control population may yield a bimodal distribution under treatment. The effect of interest in this case is to detect bimodality. A simple solution is to use the shorth functional: for any data point  $x$  and any chosen level  $\alpha$  find the length of the smallest interval containing  $x$  and covering at least a proportion  $\alpha$  of the data. You can draw this length as a function of  $x$  for various levels  $\alpha$  since the curves are not overlapping. So you do not run into the bandwidth problem of histograms or kernel density estimators. We will use this idea in the next chapter.

**ToDo:**  
add minimal guaranteed effect

### 3.7 R Complements

s:3.4



### 3.8 Statistical Summary

As a leading example in this chapter we use the comparison of treatments. In simple cases, the samples under the treatments to be compared differ only by a shift in mean. In this case the problems can be reduced to the approaches of chapter 2. In this reduced case the distributions coincide when centred at the mean. For the general case which we touched here the simplification does not work. An important example is the analysis of therapy studies. If a treatment has a homogenous effect, we can analyse it using the methods of chapter 2. But often under treatment we see a split into “responders” and “non responders”. This goes beyond the models sketched in chapter 2 and requires more general approaches mentioned in this chapter 3.

We restricted our discussion to the comparison of two samples. Practical work often leads to other problems. For example, a typical question is to compare a new treatment with a known reference treatment, where only sample observations are available for the new treatment, but extensive prior information is available about the reference treatment. Or a reference treatment is to be compared with a series of alternative treatments. These problems go beyond the scope of our introduction. Here we can only refer to advanced literature, for example [17].

### 3.9 Literature and Additional References

- [31] Venables, W.N.; Ripley, B.D. (2002): Modern Applied Statistics with S. Heidelberg: Springer
- [30] Venables, W.N.; Ripley, B.D. (2000): S Programming. Heidelberg: Springer
- [17] Miller, R. G. (1981): Simultaneous Statistical Inference. Heidelberg: Springer

```
Generated by Sweave from:
Source: /u/math/j40/cvsroot/lectures/src/SIntro/Rintro/Rnw/S03vergl.Rnw.tex,v
Revision: 1.32
Date: 2008/06/18 16:48:43
name:
Author: j40
$Source: /u/math/j40/cvsroot/lectures/src/SIntro/Rintro/Rnw/S03vergl.Rnw.tex,v $
$Revision: 1.32 $
$Date: 2008/06/18 16:48:43 $
$name: $
$Author: j40 $
```



---

## CHAPTER 4

# Dimensions 1, 2, 3, ..., $\infty$

---

ch:04

*Private note:*

*What are the essential “need to know” concepts and ideas about multivariate statistics?*

*At current critically missing: Stein shrinkage and dimension reduction.*

*Should principal component regression be included?*

*Should sliced inverse regression be included?*

*Rescaling/standardisation for principal components should be made explicit.*

In the previous chapters, we used numerical and graphical tools in parallel. The general idea is that there is a strict correspondence. A graphical tool is only worth as much as the hard statistics that is supporting it. If we go to higher dimensions, we have to review this correspondence.

Part comes from perception. It is doubtful that humans are made for handling numbers. Even if we are numerate, most of us barely can go beyond one dimension as far as numbers are concerned. We can display numbers in a table, generally with a matrix layout. This gives us two dimensions for the position, and one for the value. But if you make any experiments, you will notice that few are able to exploit this three dimensional structure. Multiway tables are even more of a challenge to read. Numbers are for the one dimensional people.

Graphical representations are more accessible in a higher dimensional setting than numerical representations, at least if the graphical structures correspond to those familiar from natural phenomena. Most people are able to see evolving three dimensional structures. So four dimensions (three for the space, one for time) seem to be tractable, and most people are capable to trace one or two qualitative features (such as colour or texture) along with the basic geometry, giving a perceptual space of four to six dimensions to start with.

While it may be interesting to investigate quantitative representations in higher dimensions, we will focus on graphical aspects in this chapter.

Graphics is related to geometry. But in a statistical context, they are not the same. Each variable may have its own scale, and these scales may have dependencies. Finding an appropriate geometry is a topic on its own. A first step in this direction is to find the proper scale of variables. In the following, we will often assume that a proper scaling has been done. This leaves a large open space which is not covered here.

To give an impression of the dragons which are lurking out there, we just mention one problem which is not covered here: Assume you have a data set with  $n$  observations  $X_i, i = 1, \dots, n$

**ToDo:**

ch.4: major rewrite

each recorded formally with a  $p$  dimensional variate. Assume that you have a measure of distance  $d(X_i, X_j,)$  for any two data points. Find a geometric representation of the data in a  $q$  dimensional variate, for example a planar representation in two-space. You will find approaches to this problem under the key word ***multidimensional scaling***. A whole branch of statistics is related to this kind of problems. We will only touch it and assume that a proper scaling has been done.

## 4.1 R Complements

ch04-lattice

In this chapter we begin with complements on R in order to concentrate on statistical questions for the remainder without disrupting the discussion by programming details. We have a look on the graphical possibilities that are at our disposition. The basic graphics model of R is oriented to possibilities of a plotter as an output device. The graphics follow the possibilities available when drawing with a pen. Besides the one and two dimensional possibilities which we have seen so far, there are possibilities to display a function which is defined over a grid. Basically three functions are available.

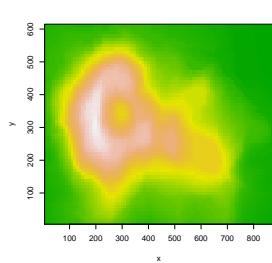
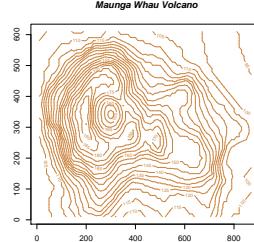
<b><i>3d Graphics</i></b>	
<i>image()</i>	gives the values of a variable $z$ in grey levels or colour coding.
<i>contour()</i>	gives the contours of a variable $z$ .
<i>persp()</i>	gives a perspective plot of a variable $z$ .

ch04-SframeVolcano *image()* and *contour()* can also be used to give an overlay on other plots.

**Example 4.1: 3d Surface Displays Using Base Graphics**

*Input*

```
#oldpar <- par(mfrow=c(1,3))
x <- 10*(1:nrow(volcano))
y <- 10*(1:ncol(volcano))
image(x, y, volcano, col = terrain.colors(100), axes = FALSE)
axis(1, at = seq(100, 800, by = 100))
axis(2, at = seq(100, 600, by = 100))
box()
title(main = "Maunga Whau Volcano", font.main = 4)
contour(x, y, volcano, levels = seq(90, 200, by = 5),
        col = "peru", main = "Maunga Whau Volcano", font.main = 4)
z <- 2 * volcano          # Exaggerate the relief
x <- 10 * (1:nrow(z))    # 10 meter spacing (S to N)
y <- 10 * (1:ncol(z))    # 10 meter spacing (E to W)
## Don't draw the grid lines : border = NA
#par(bg = "slategray")
persp(x, y, z, theta = 135, phi = 30, col = "green3", scale = FALSE,
      ltheta = -120, shade = 0.75, border = NA, box = FALSE)
```

*image()**contour()**persp()*

The basic graphics system is easy to use, but limited in possibilities. A newer graphics system conceptually works with objects and a camera model. The graphic objects can be combined and post-processed. The display takes part in a separate step. Simple 2d graphs can be post-processed. For a 3d display distance, point of view and focal length can be chosen like when using a camera. The object oriented graphics system consists of a library *grid* with the elementary operations required, and a higher level library *lattice* which gives new implementation of the displays known from the basic graphics and adds additional displays.

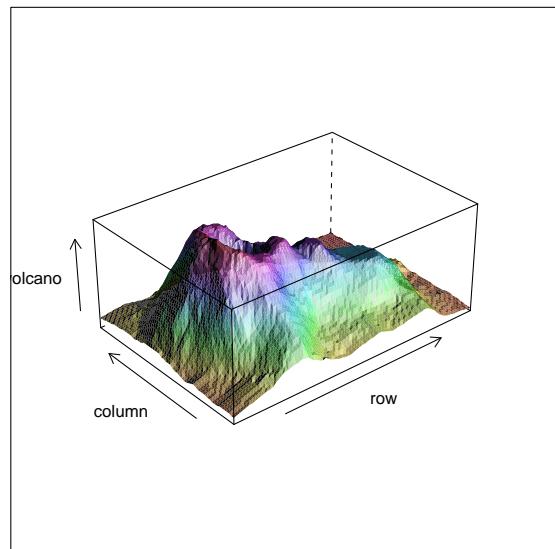
<i>3d Graphics</i>	
<i>cloud()</i>	Generic lattice function to draw 3d scatterplots.
<i>wireframe()</i>	Generic lattice function to draw 3d surfaces.

In the basic graphics system, the functions usually provide a graphical output, and the internal information must be accessed explicitly. In the lattice system, the functions usually return lattice objects. Graphical output must be requested explicitly. For the output of lattice objects the function `print()` is used.

**Example 4.2: 3d Surface Display Using Lattice Graphics**

*Input*

```
library(lattice)
## volcano ## 87 x 61 matrix
print(wireframe(volcano, shade = TRUE,
                aspect = c(61/87, 0.4),
                light.source = c(10,0,10)))
```



The basic graphics system and lattice graphics are separate graphics systems. Unfortunately they use different notations for comparable functions, and comparable displays have different representations. A small translation aid is given in table 4.5. Some convenience functions to combine both graphics systems are provided in library `gridBase`. An extensive introduction to both graphics systems is [18].

In a wide range of scientific visualisations, OpenGL is used as a common standard. OpenGL functions are accessible in R using the library `rgl`. There are however certain differences between common requirements for graphics, and the specific requirements of statistical graphics. As far as the representation of functions is concerned, statistical graphics is comparable with the requirements usual in analysis. The small difference is that functions in statistics are often piecewise constant or only piece-wise continuous, while for example in analysis continuous or even differentiable functions are the rule rather than the exception. When it comes to displaying data, the situation changes drastically. Usually, statistical data are discrete. Smoothness properties

Basic Graphics		Lattice
<code>barplot()</code>	bar chart	<code>barchart()</code>
<code>boxplot()</code>	box and whisker plot	<code>bwplot()</code>
	3 dimensional scatterplot	<code>cloud()</code>
<code>contour</code>	contour plot	<code>contourplot()</code>
<code>coplot</code>	conditional scatterplots	<code>contourplot()</code>
<code>curve(density())</code>	density estimator	<code>densityplot()</code>
<code>dotchart()</code>	dot plot	<code>dotplot()</code>
<code>hist()</code>	dot plot	<code>histogram()</code>
<code>image()</code>	colour map plots	<code>splom()</code>
	parallel coordinate plots	<code>parallel()</code>
<code>pairs()</code>	scatterplot matrices	<code>wireframe()</code>
<code>persp()</code>	three dimensional surface	<code>wireframe()</code>
<code>plot()</code>	scatterplot	<code>xyplot()</code>
<code>qqnorm()</code>	theoretical $QQ$ plot	<code>qqmath()</code>
<code>qqplot()</code>	empirical $QQ$ plot	<code>qq()</code>
<code>stripchart()</code>	one dimensional scatterplot	<code>stripplot()</code>

Table 4.5 *Basic Graphics and Lattice Graphics*

ta:ch04lattice

that simplify display of analytical data are not available for statistical data. So visualisations adapted to the needs of statistics are required.

## 4.2 Dimensions

If we move from one dimension to higher dimensions, there are challenges for theoretical considerations as well as for graphical displays. Again, linear models can serve as leading or warning examples.

The challenge may result from serious problems. Even under regularity conditions, the distribution on higher dimensional spaces can be complex beyond limits. The classification identification problems for functions and spaces from analysis and geometry give a foretaste of what has to be coped with for the analysis of probability distributions.

Besides, there are home made problems which are generated by choices made in a first step.

An example of home made problems can be illustrated for linear models. The interpretation of the Gauss-Markov estimator as a linear projection shows that seemingly only the coefficients for single regressors are estimated. In fact, a vector in the vector space spanned by the regressors is estimated; the attribution to individual regressors is mere linear algebra. The allocation to individual components does not depend on the influence of the individual regressor, but on

the joint geometry of all the regressors. Only if the regressors form an orthogonal basis there is a direct interpretation of the coefficients. For example if we duplicate the list of regressors in a linear model, the space does not change. The calculation in coordinates become more complicated because the regressors do not form a basis now, but from an abstract point of view, the information is unchanged. However if there is no exact duplicate, but minute deviations (by minimal “error”, rounding, transformations), the situation changes drastically. With respect to the Gauss-Markov estimator, only the space spanned by the regressors is relevant, and even minimal changes in the duplicate copy can duplicate the dimension of this space. This is a example of a home made problem.

This and other examples are reason enough to analyse the relation between the variables more precisely. In the case of regression for example this does not only refer to the relation between response and regressors, but, as illustrated by the recent example, the relations between the regressors.

To keep the relation to the regression problem and resort to experiences in this area, we use a formal imbedding of the regression problem in a more general framework. In regression, we had a singled out variable, the response. The distribution of this variable was to be modelled as a function of the remaining variables, the regressors  $X$ . Now we comprise response and regressor to a joined data vector  $Z = (Y; X)$  and will discuss the joined distribution of  $Z$ . We retrieve the regression problem in this more general framework: for the regression problem we were looking for an estimator for the mean value function  $m$  in the model

$$Y = m(X) + \varepsilon.$$

In the general framework we consider the joint distribution of  $X$  and  $Y$ . The regression model becomes

$$Y = E(Y|X) + \varepsilon$$

and for now we have the identification  $m(X) = E(Y|X)$ .

If we are indeed interested in the regression model, additional work has to be done. An estimation of the conditional expected value  $E(Y|X)$  is not the same as the estimation of a regression function  $m(X)$ . In the regression problem we have made no assumptions about the distribution of  $X$ . To infer from  $E(Y|X)$  (or a estimator of this) to  $m(X)$  we must establish that the estimation is independent of the distribution assumptions about  $X$ . For present purposes the distinction is irrelevant. We can allow ourselves a temporary ignorance.

So the general frame in this chapter is: we analyse data  $(Z_i)_{i=1,\dots,n}$ , where the single observations take values in  $\mathbb{R}^q$ . The regression model is contained in this setting as  $Z_i = (X_i, Y_i)$ .

If we have essentially linear structures, we often can analyse higher dimensional structures with methods that have been developed for one dimensional models. We may have to modify these methods, or use them stepwise taking residuals. But they may help us to detect the essential structures. They fail, however, if higher dimensionality combines with non linearity. In that case, more specific methods are needed.

### 4.3 Selections

Originally a selection indicated a selection of observations. For the graphical display, a selection is associated with a specification of display attributes (such as colour, plot character, line width). All variable values that belong to observations in the selection are displayed with these

Name	Variable	Unit, Remarks
rw	relative body weight	
fpg	plasma glucose (after overnight fast)	[mg/100 ml]
ga	glucose level integrated over about 3 hours tolerance test	[mg/100 ml × h]
ina	insulin level integrated over about 3 hours tolerance test	[μU/100 ml × h]
sspg	plasma glucose (steady state)	[mg/100 ml]
cc	classification	chemical, normal, open diabetes

Table 4.6 *Diabetes data set: variables***diabvars**

attributes. This allows to follow links by selection in various plots (“*linking*”). With linking, selections can help to detect structures in linked plots.

In practical data analysis, selections are varied dynamically (“*brushing*”), to group corresponding observations. Brushing with linked plots is an important tool in interactive data analysis. In statistical terms, selections are used for model selection. They correspond to local models: instead of using a global, possibly rather complex for the data, for each selection a (hopefully simpler) model is used which only applies to the data in this selection.

Unfortunately linking is not supported directly by R. We ourselves have to ensure that selections are displayed using their accordant attributes. Moreover the syntactic representation of selections is not unified. In function calls, selections can be represented as a *selection* parameter, or by a *group* variable, or as a condition in a formula expression. So in each case we have to resume to ad-hoc solutions.

R is mainly restricted to static selections. Brushing is only possible in a rudimentary form.

Selections are illustrated in context in the following sections.

#### 4.4 Projections

As a first example we look at a data set published in a comparison of different kind of diabetes ([freeman1979diab](#)). The data set is provided for example in `library(locfit) as data(chemdiab)`. The variables refer to laboratory values on glucose metabolism (see table 4.6).

---

*Input*

```
library(locfit)
data(chemdiab)
```

---

We get a first summary with

---

*Input*

---

```
summary(chemdiab)
```

---

		Output	
rw	fpg	ga	
Min.	:0.7100	Min. : 70.0	Min. : 269.0
1st Qu.	:0.8800	1st Qu.: 90.0	1st Qu.: 352.0
Median	:0.9800	Median : 97.0	Median : 413.0
Mean	:0.9773	Mean :122.0	Mean : 543.6
3rd Qu.	:1.0800	3rd Qu.:112.0	3rd Qu.: 558.0
Max.	:1.2000	Max. :353.0	Max. :1568.0
		ina                    sspg	cc
Min.	: 10.0	Min. : 29.0	Chemical_Diabetic:36
1st Qu.	:118.0	1st Qu.:100.0	Normal :76
Median	:156.0	Median :159.0	Overt_Diabetic :33
Mean	:186.1	Mean :184.2	
3rd Qu.	:221.0	3rd Qu.:257.0	
Max.	:748.0	Max. :480.0	

---

As in the original publication we omit the relative weight from our discussion. The chemical classification *cc* is derived from the metabolic data. So it does not contain any information on its own. For orientation, we use it as classification marker, that is we use the selections *cc* = *Chemical\_Diabetic*, *Normal*, *Overt\_Diabetic*. Essentially the data set is four dimensional with the variables *fpg*, *ga*, *ina*, *sspg*.

#### 4.4.1 Marginal Distributions and Scatter Plot Matrices

We can try to analyse multi-dimensional distributions by looking at the two dimensional *marginal distributions* for all pairs of variables. The corresponding graphical display is called a *scatterplot matrix*. In R, it is implemented as function *pairs()*.

---

**help(pairs)**

---



---

pairs	<i>Scatterplot Matrices</i>
-------	-----------------------------

---

#### Description

A matrix of scatterplots is produced.

#### Usage

```
pairs(x, ...)

## S3 method for class 'formula':
pairs(formula, data = NULL, ..., subset,
```

```

na.action = stats::na.pass)

## Default S3 method:
pairs(x, labels, panel = points, ...,
      lower.panel = panel, upper.panel = panel,
      diag.panel = NULL, text.panel = textPanel,
      label.pos = 0.5 + has.diag/3,
      cex.labels = NULL, font.labels = 1,
      row1attop = TRUE, gap = 1)

```

*Arguments*

<b>x</b>	the coordinates of points given as numeric columns of a matrix or data frame. Logical and factor columns are converted to numeric in the same way that <b>data.matrix</b> does.
<b>formula</b>	a formula, such as $\sim x + y + z$ . Each term will give a separate variable in the pairs plot, so terms should be numeric vectors. (A response will be interpreted as another variable, but not treated specially, so it is confusing to use one.)
<b>data</b>	a <b>data.frame</b> (or <b>list</b> ) from which the variables in <b>formula</b> should be taken.
<b>subset</b>	an optional vector specifying a subset of observations to be used for plotting.
<b>na.action</b>	a function which indicates what should happen when the data contain NAs. The default is to pass missing values on to the panel functions, but <b>na.action = na.omit</b> will cause cases with missing values in any of the variables to be omitted entirely.
<b>labels</b>	the names of the variables.
<b>panel</b>	<b>function(x,y,...)</b> which is used to plot the contents of each panel of the display.
<b>...</b>	arguments to be passed to or from methods. Also, graphical parameters can be given as can arguments to <b>plot</b> such as <b>main.par("oma")</b> will be set appropriately unless specified.
<b>lower.panel, upper.panel</b>	separate panel functions to be used below and above the diagonal respectively.
<b>diag.panel</b>	optional <b>function(x, ...)</b> to be applied on the diagonals.
<b>text.panel</b>	optional <b>function(x, y, labels, cex, font, ...)</b> to be applied on the diagonals.
<b>label.pos</b>	y position of labels in the text panel.
<b>cex.labels, font.labels</b>	graphics parameters for the text panel.
<b>row1attop</b>	logical. Should the layout be matrix-like with row 1 at the top, or graph-like with row 1 at the bottom?
<b>gap</b>	Distance between subplots, in margin lines.

### Details

The  $ij$ th scatterplot contains  $x[, i]$  plotted against  $x[, j]$ . The scatterplot can be customised by setting panel functions to appear as something completely different. The off-diagonal panel functions are passed the appropriate columns of  $x$  as  $x$  and  $y$ : the diagonal panel function (if any) is passed a single column, and the `text.panel` function is passed a single  $(x, y)$  location and the column name.

The graphical parameters `pch` and `col` can be used to specify a vector of plotting symbols and colors to be used in the plots.

The graphical parameter `oma` will be set by `pairs.default` unless supplied as an argument. A panel function should not attempt to start a new plot, but just plot within a given coordinate system: thus `plot` and `boxplot` are not panel functions.

By default, missing values are passed to the panel functions and will often be ignored within a panel. However, for the formula method and `na.action = na.omit`, all cases which contain a missing values for any of the variables are omitted completely (including when the scales are selected).

### Author(s)

Enhancements for R 1.0.0 contributed by Dr. Jens Oehlschlaegel-Akiyoshi and R-core members.

### References

Becker, R. A., Chambers, J. M. and Wilks, A. R. (1988) *The New S Language*. Wadsworth & Brooks/Cole.

### Examples

```

pairs(iris[1:4], main = "Anderson's Iris Data -- 3 species",
      pch = 21, bg = c("red", "green3", "blue") [unclass(iris$Species)])  
  

## formula method  

pairs(~ Fertility + Education + Catholic, data = swiss,  

      subset = Education < 20, main = "Swiss data, Education < 20")  
  

pairs(USJudgeRatings)  
  

## put histograms on the diagonal
panel.hist <- function(x, ...)  

{  

  usr <- par("usr"); on.exit(par(usr))  

  par(usr = c(usr[1:2], 0, 1.5) )  

  h <- hist(x, plot = FALSE)  

  breaks <- h$breaks; nB <- length(breaks)  

  y <- h$counts; y <- y/max(y)  

  rect(breaks[-nB], 0, breaks[-1], y, col="cyan", ...)  

}  

pairs(USJudgeRatings[1:5], panel=panel.smooth,

```

```
cex = 1.5, pch = 24, bg="light blue",
diag.panel=panel.hist, cex.labels = 2, font.labels=2)

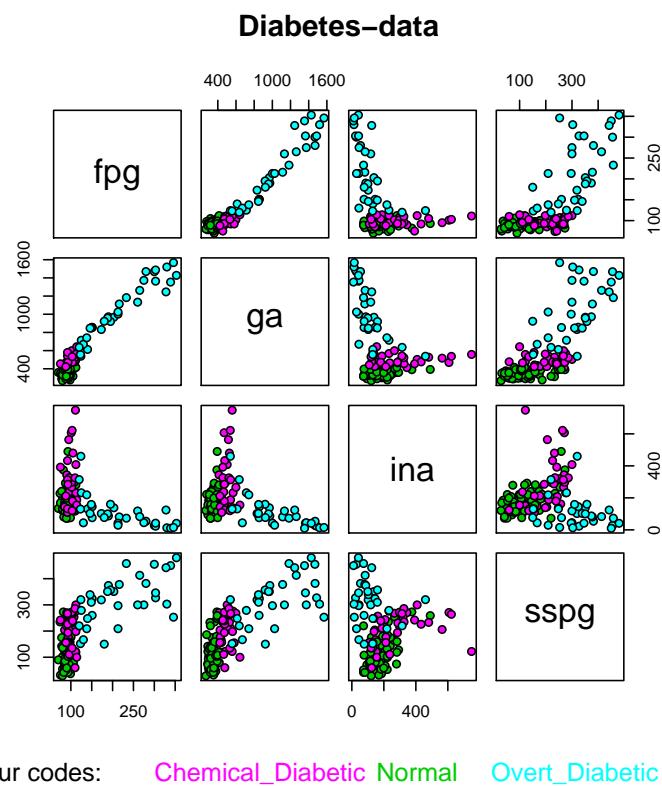
## put (absolute) correlations on the upper panels,
## with size proportional to the correlations.
panel.cor <- function(x, y, digits=2, prefix="", cex.cor)
{
  usr <- par("usr"); on.exit(par(usr))
  par(usr = c(0, 1, 0, 1))
  r <- abs(cor(x, y))
  txt <- format(c(r, 0.123456789), digits=digits)[1]
  txt <- paste(prefix, txt, sep="")
  if(missing(cex.cor)) cex.cor <- 0.8/strwidth(txt)
  text(0.5, 0.5, txt, cex = cex.cor * r)
}
pairs(USJudgeRatings, lower.panel=panel.smooth, upper.panel=panel.cor)
```

We use the chemical diabetes classes `cc` to define selections. We assign a colour to each of these selections. This is the `linking` attribute that allows us to trace relations between the plots. To provide a documentation of this linking, we have to control the graphics and modify the plots. Using the parameter `oma` we generate an outer margin for a caption.

expl:ch04-diabetes

*Example 4.3: Scatterplot Matrix for Diabetes Data*

```
Input
pairs(~fpg + ga + ina + sspg, data = chemdiab, pch = 21,
      main = "Diabetes-data",
      bg = c("magenta", "green3", "cyan")[unclass(chemdiab$cc)],
      oma = c(8, 8, 8, 8))
mtext(c("Colour codes:", levels(chemdiab$cc)),
      col = c("black", "magenta", "green3", "cyan"),
      at = c(0.1, 0.4, 0.6, 0.8), side = 1, line = 2)
```

**ToDo:**

supply  
better  
support  
~~fa:diag:4~~  
colour key

The `pairs()` function only controls the “layout” of the matrix, the selection and placement of the projections. The display in the individual plot tiles can be controlled in the function call. The defaults give the names of the variables in the diagonal tiles, and pair-wise scatterplots in the off diagonal tiles.

Exercise 4.1	
	Generate a scatterplot matrix for the diabetes data set that show a histogram of the variables in the diagonal tiles. <i>Hint:</i> See <code>help(pairs)</code> .

Certain aspects of the distribution can be read easily from the marginal distributions. Other geometrical structures are barely or not at all recognisable from the marginal distribution.

For example, there is an apparent linear relation between the glucose level after a fast `fpg` and the integrate glucose level under tolerance test. This can be seen from the two dimension marginal distributions, and it can be analysed using the methods for linear models.

This apparent relation propagates to the relations to the other variables `ina`, `sspg`. For this reason, `fpg` is not considered in detail in the original paper. What remains to analyse are the variables `ga`, `ina`, `sspg`. The three dimensional structure of this part of the data set is not easily detected from the marginal distributions.

#### 4.4.2 Projection Pursuit

Geometrical relations or stochastic dependencies that are not aligned along co-ordinate axes are in general not covered by marginal distributions. We can generalise the idea. Instead of two dimensional marginal distributions, we can use arbitrary projections. We use `library(lattice)` which implements a camera oriented graphic model.

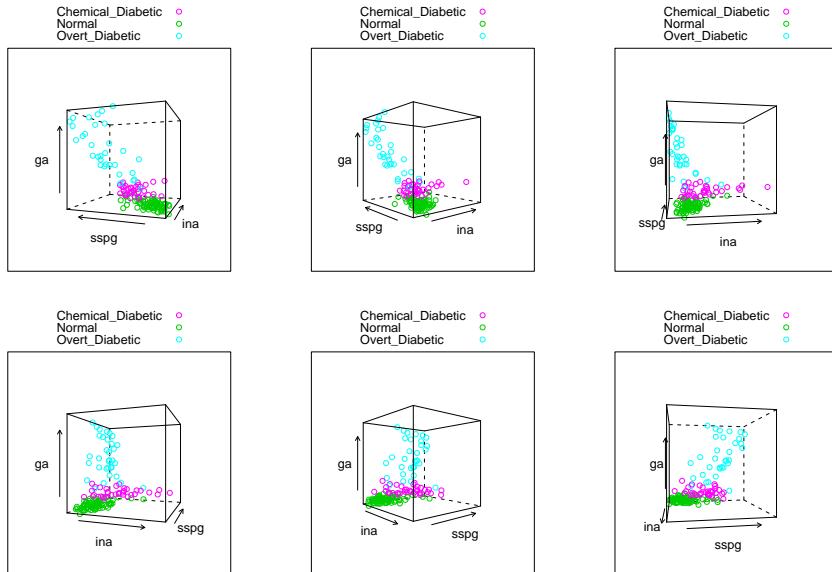
The `grid` graphics and the `lattice` package provide support for multivariate representations. `grid` forms the basis. The original R graphics system implements a model that takes pen and paper as a model. A graphic port (paper) is opened, and lines, points ore symbols are drawn. `grid` is a second graphics system which follows a camera/object model. Graphical objects in in various locations and orientations are mapped in a visual space. `lattice` builds upon `grid`. In <http://cm.bell-labs.com/cm/ms/departments/sia/project/trellis/> gives the underlying ideas for visualisation of multidimensional data that are implemented in `lattice`.

The graphics thus generated is output with `print()`. The parameter `split` allows dividing up the output area. Unfortunately linking is broken here: `cloud()` can generate a caption, but this shows the colour scale which was in effect when the graphics started, not the colour scale effectively used in the output. Once again we have to interfere with the system, this time to adjust the colour table.

*Input*

```
library("lattice")
diabcloud <- function(y, where, more = TRUE, ...) {
  print(cloud(ga ~ ina + sspg, data = chemdiab, groups = cc,
              screen = list(x = -90, y = y), distance = .4, zoom = .6,
              auto.key = TRUE, ...),
        split = c(where, 3, 2), more = more)
}
supsym <- trellis.par.get("superpose.symbol")
supsymold <- supersym
supsym$col = c("magenta", "green3", "cyan")
trellis.par.set("superpose.symbol" = supersym)
diabcloud(y = 70, where = c(1, 1))
```

```
diabcloud(y = 40, where = c(2, 1))
diabcloud(y = 10, where = c(3, 1))
diabcloud(y = -20, where = c(1, 2))
diabcloud(y = -50, where = c(2, 2))
diabcloud(y = -80, where = c(3, 2), more = FALSE)
trellis.par.set("superpose.symbol" = supsymold)
rm(diabcloud, supersymold, supersym)
```



a:ch04:1

Exercise 4.2	
	<p>Modify this example so that you get an impression of the three dimensional structure.</p> <p>What is the difference between open diabetes and chemical diabetes?</p> <p>How does the normal group compare to both diabetes groups?</p>

As with series of projections it is often not simple to identify a three dimensional structure. Animation can help. There is support for animation in *library(rggobi)*. This however requires ggobi, accessible in in <http://www.ggobi.org/>, as additional software.

What has been done here ad-hoc can be done systematically and generalised for arbitrary dimensions: for a data set in  $\mathbb{R}^q$  one looks for “interesting” projections. For this, an index is defined that should indicate how interesting a projection is. Then a search is started to maximise this index. A family of statistical methods based on this idea can be found under the key word **projection pursuit**. *ggobi* contains implementations of projection pursuit for several indices. These can be accessed using the R functions in *library(rggobi)*.

### 4.4.3 Projections for Dimensions 1, 2, 3, ... 7

Projection methods try to identify structures of lower dimension in a higher dimensional data set. The dimensions that can be identified are limited. If we project a structure with a dimension larger than the projection target, a typical projection covers the target space and no longer gives information.

How many dimensions can we capture? The graphical display in the plane gives, as lower bound, two dimensions. We can represent two dimensional structures directly using Cartesian co-ordinates in the  $xy$  plane (at least locally. The global structure might require additional dimensions.) Perception can reconstruct three dimensional structures with the help of visual depth cues (such as a shade) or from sequences of 2d images. Animation can give the impression of changing 3d scenes and we reach four dimensions.

Additional information channels, as for example colour encoding can elevate the dimension slightly, but effectively we stay with four to seven dimensions for a display.

Combination of several displays helps little beyond this limit. If we combine several displays, for example generalising the idea of a scatterplot matrix, we lose the capability to generate more complex structures by our perception. Instead we have to work out the more complex structures by active comparison. The ability to perform simultaneous comparison of structures is limited. The same, is the number of displays which can be presented simultaneous on a page medium such as a screen or paper is limited.

Four dimensions is easy. Going beyond seven is hard.

For illustration for an example in seven dimensions, we use a classical experiment on colour perception. An observer (usually a person well trained in working with colours) is presented a target colour, and control over three light sources (red, green blue). The task is to adjust the controllable light sources to match the target colour. The classical experiment was performed with 34 target colours varying in frequency in steps of 10nm, with 20 observers.

To start with, we can view a single observation as a four dimensional datum (target colour, intensities for red, green, blue). We can take the three marginals of the target colour by any of the intensities. This gives us the colour matching curves which are found in many text books on colour perception. (See [fig:04color](#) top row). We use colour encoding by target colour here as an additional cue. Several notes are necessary. First, the classical experiments have been made with a black background, at that time. Second, the colour cues are a rough approximation, not the exact target colour. Third, not all colours can be mixed from the primaries in this experiment. The observer had the possibility to set the control to negative values which would have the effect to change the target colour by blending. Fourth, as in all real data, there are measurement errors. One observation is off.

The basic data set has four dimensional observations. The target colour is already encoded as colour cue. We can combine the three response dimensions in 3d scatterplot. An example is in figure [fig:04color](#), bottom row middle. It tells us that essentially we have a one dimensional data structure, imbedded non linearly in 3d space.

From the experience with linear models we know that the data scatterplot only shows part of the information. The other information is contained in the residuals. Taking residuals from one dimensional non linear regression for the three components (figure [fig:04color](#), middle row) gives us another three dimensional view on the information contained in the data. So the complete data structure has seven dimensions, with one dimension (the target colour) is used for linking and represented by colour. Again using a (non-marginal) projection from this 3d structure tells us

**ToDo:**  
Does  
memory  
help to in-  
crease the  
bound?

**ToDo:**  
color R  
version

4-16

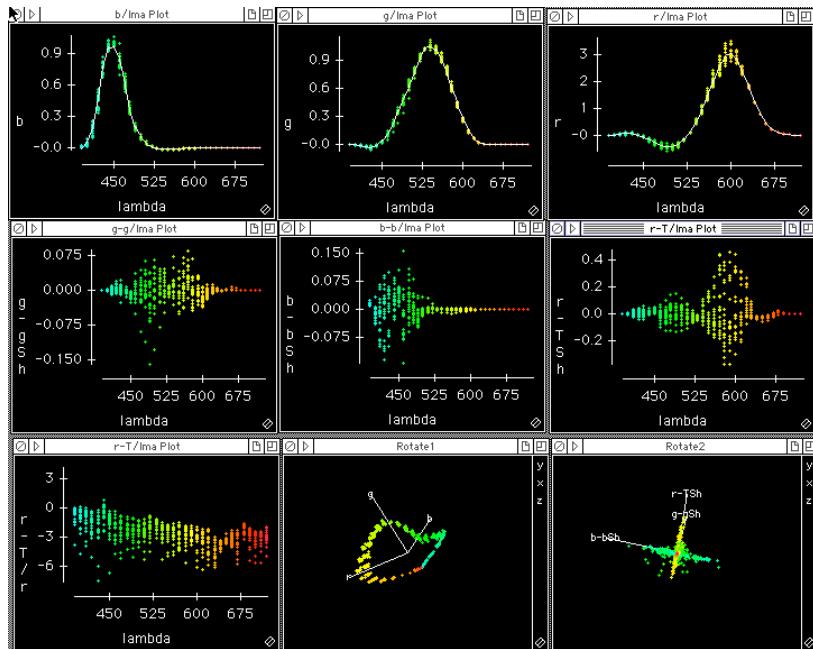
DIMENSIONS 1, 2, 3, ...,  $\infty$ 

Figure 4.1 Colour perception data. Displays prepared with Data Desk

fig:04color

that again we have a non linear structure: there is a disk, essentially in the plane characterised by red/green colours, and rod-like structure orthogonal to this, essentially characterised by blue/yellow colours.

We only used basic data analysis approaches to reveal this structure. Work on colour perception has proceeded. The conception nowadays is that there are essentially two systems in early colour perception, one combining the red/green channels, the other one based on the blue channel. We did not use this structural information. But the information we get from the residuals just fit.

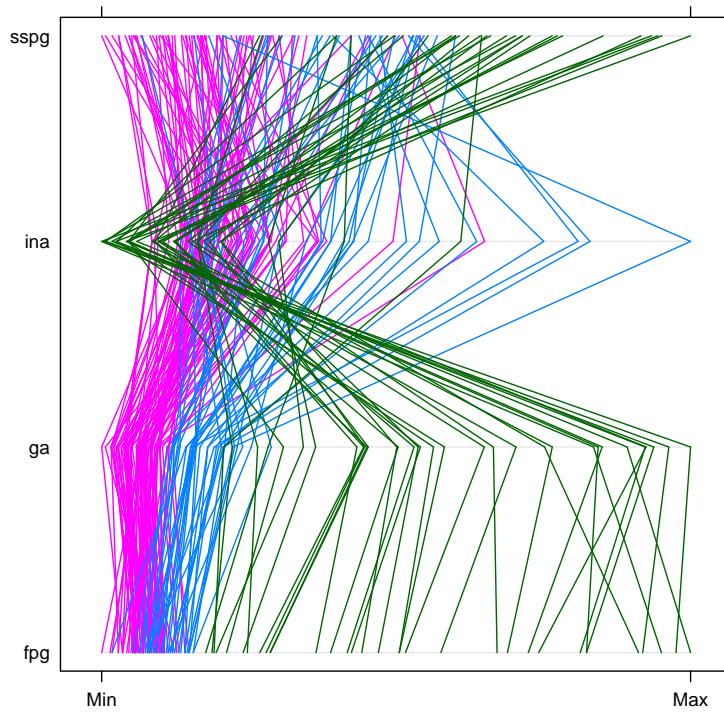
#### 4.4.4 Parallel Coordinates

The graphical display (in Cartesian co-ordinates) is initially restricted to one and two dimensional projections. But even for representation in a plane, the restriction to two dimensions is not predetermined, but is a consequence of our choice to display the data in Cartesian co-ordinates. Plot matrices break this dimension barrier by combining Cartesian co-ordinate systems.

Parallel co-ordinates use a parallel orientation of axes, not an orthogonal orientation. For frequencies with categorical variables this is a common display: (possibly stacked) bar charts use parallel co-ordinates. Function `parallel()` in `library(lattice)` implements parallel co-ordinates for quantitative variables as well. The marks on these axes that correspond to one case are joined by a zigzag line. This form of parallel co-ordinates was suggested by A. Inselberg ([\[13\]](#)).

*Input*

```
library("lattice")
print(parallel(chemdiab[2:5], groups = chemdiab$cc))
```



The information is the same as that in previous plots of this data set. Using a different display gives a new look on the structure in this data set.

a:ch04:2

**ToDo:**  
add legend, control colour

<b>Exercise 4.3</b>	
	For the <code>chemdiab</code> data set, prepare a (written!) report about the relation between the variables that you can recognise in the parallel co-ordinate plot.
	Instead of using <code>chemdiab[2:5]</code> you can specify the variables explicitly as <code>chemdiab[c(2, 3, 4, 5)]</code> . This gives you control over the order of the variables. Compare two different sequences of the variables and note (in writing!) your observations. Which sequence of variables gives the simples display? Which relations between the variables are visible in both? Which relations appear only in one of the arrangements?

## 4.5 Sections, Conditional Distributions and Coplots

s:4.3

From an abstract point of view, sections are conditional distributions of the type  $P(\cdot \mid X = x)$ . But they are only reliable where the section defines a condition which has positive probability. To make the idea to restrict the view on conditional distributions applicable to data, we thicken the sections. Instead of considering conditional distributions of the type  $P(\cdot \mid X = x)$  we consider  $P(\cdot \mid \|X - x\| < \varepsilon)$ , where  $\varepsilon$  possibly can vary with  $x$ . In graphical representations of data this requires a series of plots showing only the part of the data set specified by the condition.

Statistically projections lead to marginal distributions and sections to conditional distributions. In a certain sense, sections and projections are complementary: projections show structural features of low dimension. Sections are helpful to detect structural features of low codimension. For data analysis, both can be combined. The interplay of projections and sections is discussed in [8]. Like the dimension boundaries for projections there are boundaries for the codimension when using sections. We can only catch structures of small codimension. If the codimension is too large, a typical section is empty, hence it has no information.

**To Do:**  
RB define  
codimen-  
sion

As a first tool, R provides the possibility analyse two variables **conditioned** on one or more additional variables. As a graphical display **coplot()** serves for this purpose. It is a variant of the plot matrix and shows in each tile the scatterplot of two variables, given the condition.

The coplot can be inspected for pattern. If the variables shown are stochastically independent of the conditioning variables, all plot elements show the same shape. The variables shown and the conditioning variables can then be de-coupled.

If the general shape coincides, but location and size vary, this hints to a (not necessarily linear) shift/scale relation. Additive models or variants of these can be used to model the relation between the variables shown and conditioning variables.

If the shape changes with varying condition, a major dependency structure or interaction my apply that needs more precise modelling.

**help(coplot)**

**coplot**

*Conditioning Plots*

*Description*

This function produces two variants of the **conditioning plots** discussed in the reference below.

*Usage*

```
coplot(formula, data, given.values, panel = points, rows, columns,
       show.given = TRUE, col = par("fg"), pch = par("pch"),
       bar.bg = c(num = gray(0.8), fac = gray(0.95)),
       xlab = c(x.name, paste("Given :", a.name)),
       ylab = c(y.name, paste("Given :", b.name)),
       subscripts = FALSE,
       xlabel = function(f) abbreviate(levels(f)),
       number = 6, overlap = 0.5, xlim, ylim, ...)
co.intervals(x, number = 6, overlap = 0.5)
```

*Arguments*

<b>formula</b>	a formula describing the form of conditioning plot. A formula of the form $y \sim x   a$ indicates that plots of $y$ versus $x$ should be produced conditional on the variable $a$ . A formula of the form $y \sim x   a * b$ indicates that plots of $y$ versus $x$ should be produced conditional on the two variables $a$ and $b$ .
	All three or four variables may be either numeric or factors. When $x$ or $y$ are factors, the result is almost as if <code>as.numeric()</code> was applied, whereas for factor $a$ or $b$ , the conditioning (and its graphics if <code>show.given</code> is true) are adapted.
<b>data</b>	a data frame containing values for any variables in the formula. By default the environment where <code>coplot</code> was called from is used.
<b>given.values</b>	a value or list of two values which determine how the conditioning on $a$ and $b$ is to take place. When there is no $b$ (i.e., conditioning only on $a$ ), usually this is a matrix with two columns each row of which gives an interval, to be conditioned on, but it can also be a single vector of numbers or a set of factor levels (if the variable being conditioned on is a factor). In this case (no $b$ ), the result of <code>co.intervals</code> can be used directly as <code>given.values</code> argument.
<b>panel</b>	a <code>function(x, y, col, pch, ...)</code> which gives the action to be carried out in each panel of the display. The default is <code>points</code> .
<b>rows</b>	the panels of the plot are laid out in a <code>rows</code> by <code>columns</code> array. <code>rows</code> gives the number of rows in the array.
<b>columns</b>	the number of columns in the panel layout array.
<b>show.given</b>	logical (possibly of length 2 for 2 conditioning variables): should conditioning plots be shown for the corresponding conditioning variables (default <code>TRUE</code> )
<b>col</b>	a vector of colors to be used to plot the points. If too short, the values are recycled.
<b>pch</b>	a vector of plotting symbols or characters. If too short, the values are recycled.
<b>bar.bg</b>	a named vector with components " <code>num</code> " and " <code>fac</code> " giving the background colors for the (shingle) bars, for <b>numeric</b> and <b>factor</b> conditioning variables respectively.
<b>xlab</b>	character; labels to use for the $x$ axis and the first conditioning variable. If only one label is given, it is used for the $x$ axis and the default label is used for the conditioning variable.
<b>ylab</b>	character; labels to use for the $y$ axis and any second conditioning variable.
<b>subscripts</b>	logical: if true the panel function is given an additional (third) argument <code>subscripts</code> giving the subscripts of the data passed to that panel.
<b>axlabels</b>	function for creating axis (tick) labels when $x$ or $y$ are factors.
<b>number</b>	integer; the number of conditioning intervals, for $a$ and $b$ , possibly of length 2. It is only used if the corresponding conditioning variable is not a factor.

<code>overlap</code>	numeric < 1; the fraction of overlap of the conditioning variables, possibly of length 2 for <code>x</code> and <code>y</code> direction. When <code>overlap</code> < 0, there will be <i>gaps</i> between the data slices.
<code>xlim</code>	the range for the <code>x</code> axis.
<code>ylim</code>	the range for the <code>y</code> axis.
<code>...</code>	additional arguments to the panel function.
<code>x</code>	a numeric vector.

### Details

In the case of a single conditioning variable `a`, when both `rows` and `columns` are unspecified, a ‘close to square’ layout is chosen with `columns >= rows`.

In the case of multiple `rows`, the *order* of the panel plots is from the bottom and from the left (corresponding to increasing `a`, typically).

A panel function should not attempt to start a new plot, but just plot within a given coordinate system: thus `plot` and `boxplot` are not panel functions.

The rendering of arguments `xlab` and `ylab` is not controlled by `par` arguments `cex.lab` and `font.lab` even though they are plotted by `mtext` rather than `title`.

### Value

`co.intervals(., number, .)` returns a (`number`  $\times$  2) matrix, say `ci`, where `ci[k,]` is the range of `x` values for the `k`-th interval.

### References

- Chambers, J. M. (1992) *Data for models*. Chapter 3 of *Statistical Models in S* eds J. M. Chambers and T. J. Hastie, Wadsworth & Brooks/Cole.  
 Cleveland, W. S. (1993) *Visualizing Data*. New Jersey: Summit Press.

### See Also

`pairs`, `panel.smooth`, `points`.

### Examples

```
## Tonga Trench Earthquakes
coplot(lat ~ long | depth, data = quakes)
given.depth <- co.intervals(quakes$depth, number=4, overlap=.1)
coplot(lat ~ long | depth, data = quakes, given.v=given.depth, rows=1)

## Conditioning on 2 variables:
ll.dm <- lat ~ long | depth * mag
coplot(ll.dm, data = quakes)
coplot(ll.dm, data = quakes, number=c(4,7), show.given=c(TRUE,FALSE))
coplot(ll.dm, data = quakes, number=c(3,7),
       overlap=c(-.5,.1)) # negative overlap DROPS values
```

```

## given two factors
Index <- seq(length=nrow(warpbreaks)) # to get nicer default labels
coplot(breaks ~ Index | wool * tension, data = warpbreaks,
       show.given = 0:1)
coplot(breaks ~ Index | wool * tension, data = warpbreaks,
       col = "red", bg = "pink", pch = 21,
       bar.bg = c(fac = "light blue"))

## Example with empty panels:
with(data.frame(state.x77), {
  coplot(Life.Exp ~ Income | Illiteracy * state.region, number = 3,
         panel = function(x, y, ...) panel.smooth(x, y, span = .8, ...))
  ## y ~ factor -- not really sensical, but 'show off':
  coplot(Life.Exp ~ state.region | Income * state.division,
         panel = panel.smooth)
})

```

---

We illustrate coplots with the “Quakes” data set. This data set gives the geographical longitude and latitude of a series of earth quakes near the Fiji island, together with the depth of the seismic centre. We use the geographical longitude and latitude as variables for marginal projection, and the depth as covariate defining the sections.

We re-code the depth, so that in graphical representations large depth points downwards.

**To Do:**  
or change  
plot co-  
ordinates?

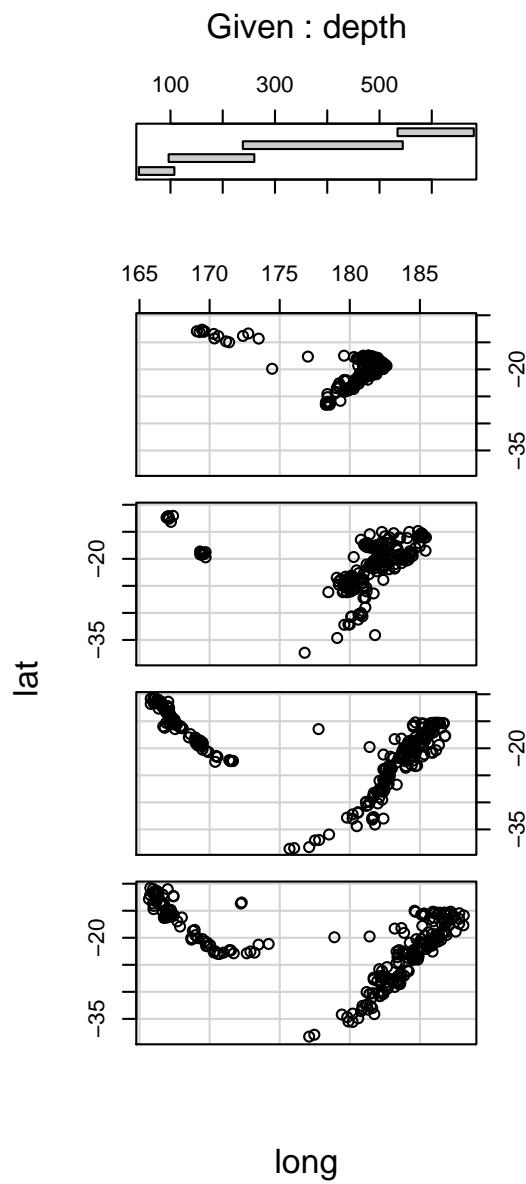
*Input*

---

```

quakes$depth <- -quakes$depth
given.depth <- co.intervals(quakes$depth, number = 4, overlap = .1)
coplot(lat ~ long | depth, data = quakes, given.values = given.depth, columns = 1)

```



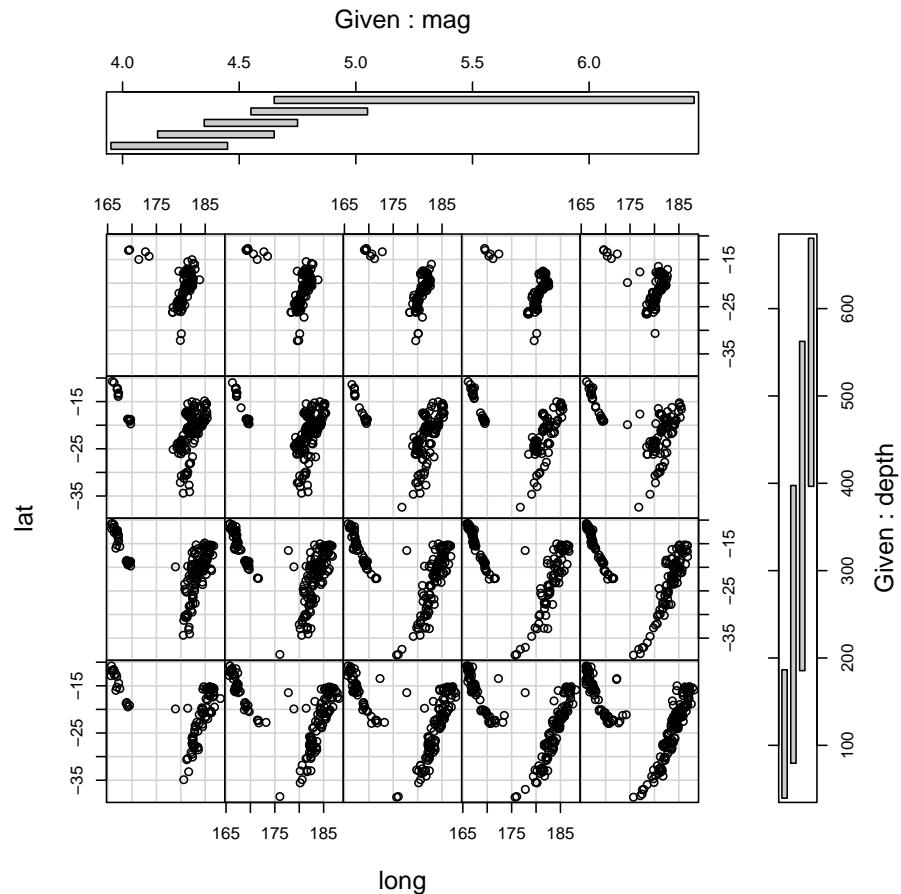
Analogue for two covariables, the depth and the magnitude of the earth quake.

---

*Input*

```
coplot(lat ~ long | mag* depth , data = quakes, number = c(5, 4))
```

---



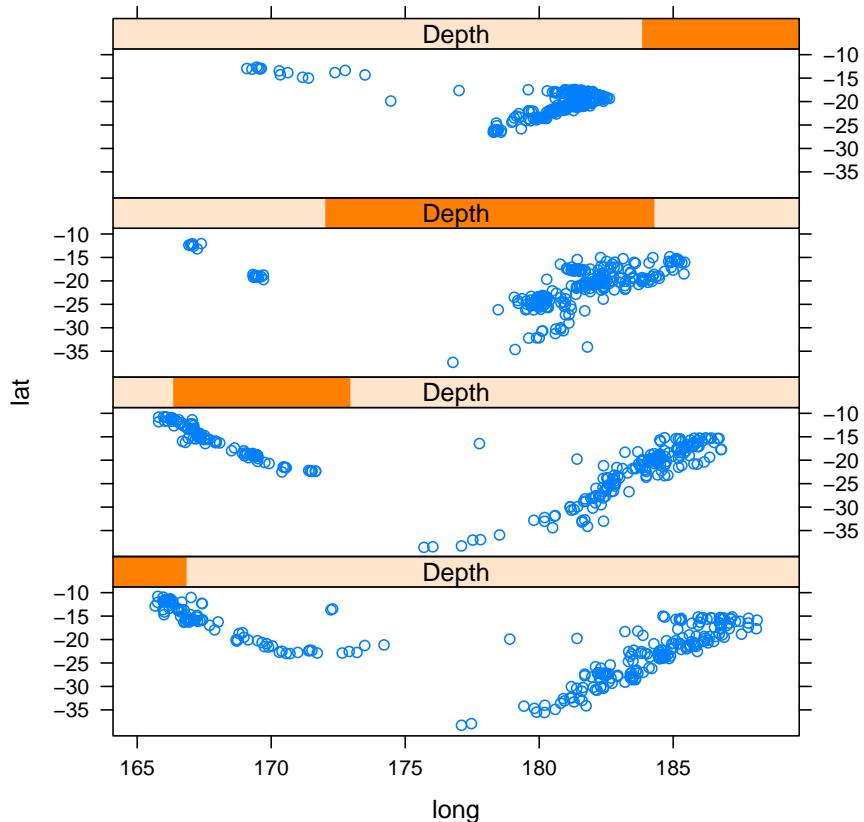
a:ch04:5

Exercise 4.4	
	Analyse the “quakes”-data set. Summarise your results in a report. Try to specify a formal model.
	How is the geographic position related to the depth?
	Can you identify relations between depth and magnitude of the earth quake? (You may have to choose a different model formula in <code>coplot()</code> .)

The idea of the coplots is generalised in the trellis displays (see [7]). Trellis displays are implemented in R in `library("lattice")`. Input [clev93vd](#)

```
library("lattice")
Depth <- equal.count(quakes$depth, number = 4, overlap = .1)
print(xyplot(lat ~ long | Depth , data = quakes, columns = 1, layout = c(1, 4)))
```

**To Do:**  
Comment  
on trellis



**ToDo:**

add magnitude

**ToDo:** Ref to curse of dimension

**ToDo:**

Variance bias dilemma

**ToDo:**

Stein shrinkage

## 4.6 Transformations and Dimension Reduction

Variables often come in a scale that is suggested by the measurement process or by conventions of the trade. This does not necessarily correspond to what would be suggested by the subject matter or what is best used for statistical modelling. This scale has some arbitrariness:

- For an acoustic perception experiment for example the strength of the stimulus can be described by energy or by acoustic pressure [phone]. One scale is transformed to the other by a logarithm. The Weber-Fechner law in psychology says that for the human perception the (logarithmic) phone scale is appropriate.
- Fuel consumption is reported in the USA in miles per gallon, in Europe as litres per 100 km. Up to conversion constants, one variable is the inverse of the other. When we want to model fuel consumption, the description in litres per 100 km seems to lead to simpler statistical models in terms of weight, power and other parameters describing the car. Analysis in terms of miles per gallon can become rather complicated.

The choice of an appropriate scale for the variables can be a critical step in the analysis. Sometimes it is helpful to first introduce transformations and constructed variables, and then to reduce dimensions and to determine the effective variables in a second step.

Co-ordinate systems are not pre-set canonically. This already applies to univariate problems. For univariate problems, transformations of co-ordinate systems are relatively simple. Modelling of the error distribution on the one hand and transformation the data to a standard distribution on the other are in a certain way exchangeable. In higher dimension situations, suitable transformation families sometimes are not available or not accessible, and the structure of a problem may critically depend on the choice of suitable co-ordinates. A choice of the co-ordinate representation based on the subject matter often is preferable to an automatic selection.

We illustrate this with Anderson's Iris-data set. The data set has five dimensions: four quantitative variables (length and width of petal sepal of iris blossoms) and a categorical variable (the species: iris setosa canadensis, iris versicolor, iris virginica)<sup>1</sup>. The question is to find a classification of the species by the four quantitative variables .

Table 4.12 *Iris species*



The structure is similar to that of the diabetes data set *chemdiab*. The classification by *iris\$Species* however is here a given (external) classification, in contrast to the derived classification *chemdiab\$cc*. While in the case of *chemdiab* a general description was looked for, the aim is now to find a classification rule that classifies *iris\$Species* by means of the other variables.

In this example, the species define the selection groups which are of interest.

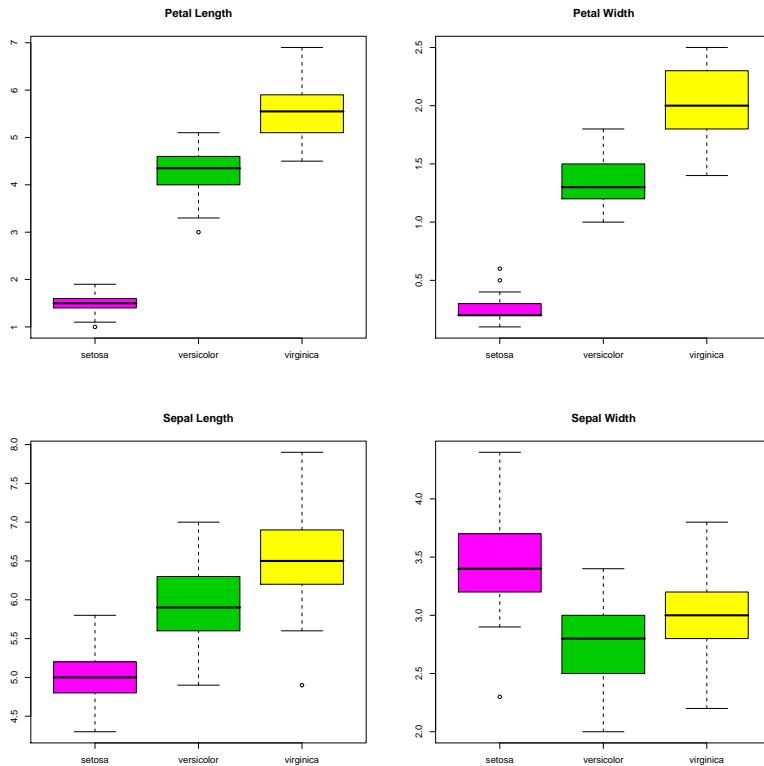
To get a first impression we can view the four quantitative variables by species. The standard conventions of R make this cumbersome. Since species is a categorical variable, this cause R to switch from a scatterplot display to box&whisker plots when using *plot()*.

---

```
oldpar <- par(mfrow = c(2, 2))
plot(iris$Species, iris$Petal.Length,
     ylab = '', main = 'Petal Length', col = c("magenta", "green3", "yellow"))
plot(iris$Species, iris$Petal.Width,
     ylab = '', main = 'Petal Width', col = c("magenta", "green3", "yellow"))
```

<sup>1</sup> Photos: Iris setosa: C. Hensler, The Rock Garden; Iris versicolor and Iris virginica: D. Kramb; The Species Iris Group of North America. With kind permission

```
plot(iris$Species, iris$Sepal.Length,
     ylab = '', main = 'Sepal Length', col = c("magenta", "green3", "yellow"))
plot(iris$Species, iris$Sepal.Width,
     ylab = '', main = 'Sepal Width', col = c("magenta", "green3", "yellow"))
par(oldpar)
```



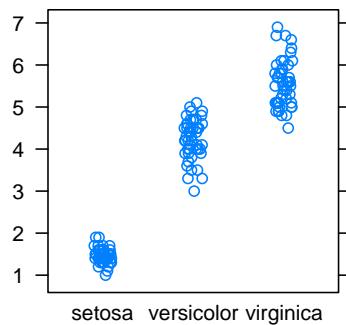
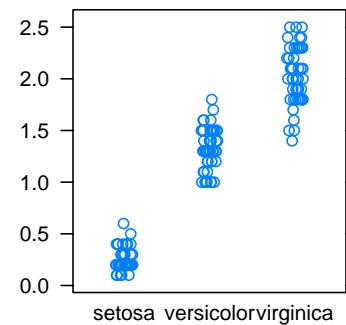
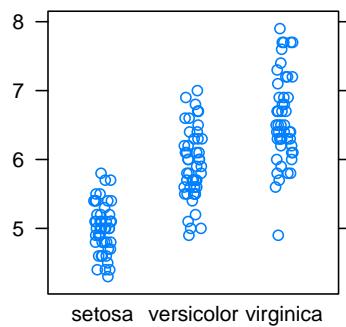
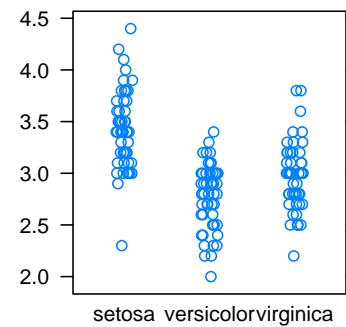
We could modify the R functions to show a scatterplot of the individual variables by group. Instead we resume to `grid` and `lattice` and use the function `stripplot()`. Since given the measurement precision in this experiment the same value is recorded in multiple cases, we jitter the data to present them separately.

---

*Input*

```
library("lattice")
print(stripplot(Petal.Length ~ Species, data = iris,
                jitter = TRUE, ylab = '', main = 'Petal Length'), split = c(1, 1, 2, 2), more = TRUE)
print(stripplot(Petal.Width ~ Species, data = iris,
                jitter = TRUE, ylab = '', main = 'Petal Width'), split = c(2, 1, 2, 2), more = TRUE)
print(stripplot(Sepal.Length ~ Species, data = iris,
                jitter = TRUE, ylab = '', main = 'Sepal Length'), split = c(1, 2, 2, 2), more = TRUE)
print(stripplot(Sepal.Width ~ Species, data = iris,
                jitter = TRUE, ylab = '', main = 'Sepal Width'), split = c(2, 2, 2, 2))
```

---

**Petal Length****Petal Width****Sepal Length****Sepal Width**

The one dimensional marginal distributions do not suffice to separate the three groups.

A two dimensional display does not help much further.

**To Do:**  
add color  
**To Do:**  
Remark  
on small  
setosa

a:ch04:11

Exercise 4.5	
	Use the methods from section 4.4 and 4.5 to inspect the data set. Can you see classification rules which give a classification of the three species to a large extent?

Using formal methods such as linear discriminant analysis (for example `lda()` in library (*MASS*)), a classification based on the morphometric variables can be found. Separating the species is not trivial.

The original variables only represent an aspect of the data that is technically most accessible. From a biological point of view, a different parametrisation would be chosen. The variables

reflect size and form of the leaves. A first approximation would be

$$\text{area} = \text{length} \cdot \text{width} \quad (4.1)$$

$$\text{aspectratio} = \text{length}/\text{width}. \quad (4.2)$$

**ToDo:**

check iris  
inverse  
ratio

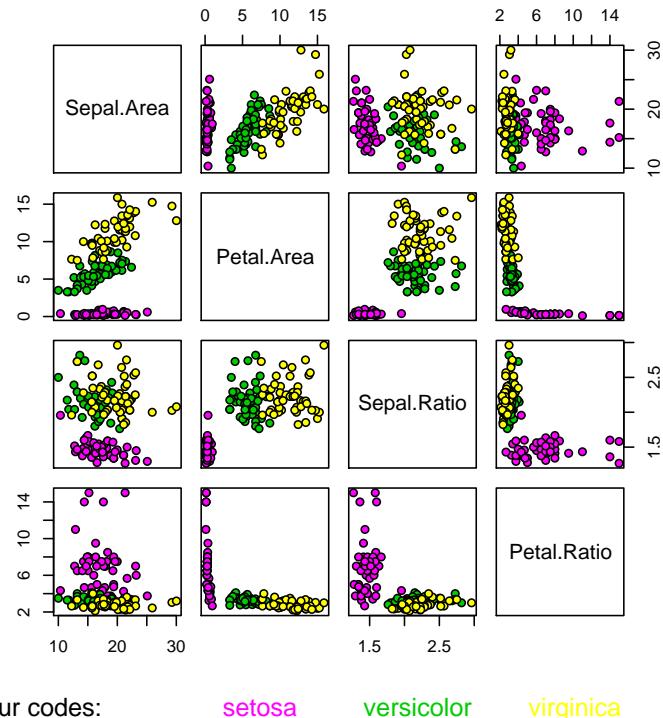
This gives the display

---

```
iris$Sepal.Area <- iris$Sepal.Length*iris$Sepal.Width
iris$Petal.Area <- iris$Petal.Length*iris$Petal.Width
iris$Sepal.Ratio <- iris$Sepal.Length/iris$Sepal.Width
iris$Petal.Ratio <- iris$Petal.Length/iris$Petal.Width
pairs(iris[6:9], main = "Anderson's Iris Data -- 3 species",
      pch = 21,
      bg = c("magenta", "green3", "yellow")[unclass(iris$Species)],
      oma = c(8, 8, 8, 8))
mtext(c("Colour codes:", levels(iris$Species)),
      col = c("black", "magenta", "green3", "yellow"),
      at = c(0.1, 0.4, 0.6, 0.8),
      side = 1, line = 2)
```

---

Anderson's Iris Data -- 3 species



Colour codes:

setosa

versicolor

virginica

**ToDo:** use  
attach

In the marginal distribution the species are separated almost completely, with two borderline cases. In these more biological co-ordinates we see that for classification the area and aspect ratio of the petal is sufficient. Any proposal for a better formal procedure has at first to improve on this trivial classification rule.

Even an exhaustive search, for example using projection pursuit, only catches projections, that is it only catches special linear combination of the variables. For the iris data, we first introduced new variables, the areas and aspect ratios. These are non linear transformations. Only in a second step the classifying variables are identified and the dimension is reduced drastically. In genuinely multivariate problems it is quite common that initially a dimension extension is necessary to solve the problem. Dimension reduction only makes sense if the describing variable are sufficiently complex to lead to a solution.

## 4.7 Higher Dimensions

### 4.7.1 Linear Case

If we have essentially linear structures, we can often analyse higher dimensional structures with methods that have been designed for one dimensional models. We may have to modify the methods or use them in iterations, taking residuals at each step. But nevertheless they may help us to recognise essential features.

In chapter 2<sup>ch:02</sup> we have already introduced linear models for arbitrary regressor dimensions  $p$ . However chapter 2<sup>ch:02</sup> presupposes that the model for the statistical analysis is specified off hand, that is the information from the data material does not influence the selection of the model, but only decisions within the framework of the model chosen.

In higher dimensional problems, it is quite common that the model is still to be specified. A common special case is the selection of regressors: the variables are candidates, among which a (preferably small) number of regressors are to be chosen. Do more complicated models provide an essential improvements over a simple model? Which parameters resp. which derived variables should be included in the model? The lesson from linear models is that not the value of individual parameters determine the contribution in a model, but that space determined by the parameters is the essential factor. Adapted strategies are needed here. We can start with simple models and ask whether additional parameter give an essential contribution. Additional parameters will improve the fit, but on the other hand the variance of the estimation will increase. We can as well start with a relatively complex model and ask whether we can omit parameters. This may increase the residual error, but we may gain reliability of the estimations.

For an abstract linear model, both strategies lead to a comparison of two model spaces  $\mathcal{M}_{X'} \subset \mathcal{M}_X$ . The corresponding estimators are  $\pi_{\mathcal{M}_{X'}}(Y)$  and  $\pi_{\mathcal{M}_X}(Y)$ . The relation between both becomes clear if we choose the orthogonal decomposition  $\mathcal{M}_X = \mathcal{M}_{X'} \oplus L_X := M_0, L_X := \mathcal{M}_X \ominus \mathcal{M}_{X'}$  of  $\mathcal{M}_X$ . Then  $\pi_{\mathcal{M}_X}(Y) = \pi_{\mathcal{M}_{X'}}(Y) + \pi_{L_X}(Y)$ .

### Partial Residuals and Added Variable Plots

In regression  $\mathcal{M}_{X'}$  and  $\mathcal{M}_X$  are spaces which are spanned by the vectors of regressor variables. In our situation we are interested in the special case

$$X' = \text{span}(X_1, \dots, X_{p'}); X = \text{span}(X_1, \dots, X_p)$$

**ToDo:**  
RB preferably small: why?

**ToDo:**  
formal variance bias

with  $p > p'$ . Here  $L_X$  is spanned by the vectors

$$R_{p'+1} = X_{p'+1} - \pi_{\mathcal{M}'_X}(X_{p'+1}), \dots, R_p = X_p - \pi_{\mathcal{M}'_X}(X_p).$$

If we (formally) do a linear regression of the additional regressors by the regressors already contained in  $X'$ , the resulting residuals generate  $L_X$ . An additional regression of  $Y$  by these residuals yields the term  $\pi_{L_X'}(Y)$  that describes the difference between the models. By construction we know that  $\pi_{\mathcal{M}'_X}(Y)$  is orthogonal to  $L_X$ . So this part is mapped to zero by the second regression. We can as well eliminate this part and restrict us to the regression of  $Y' = Y - \pi_{\mathcal{M}'_X}(Y)$  after  $R_{p'+1}, \dots, R_p$ .

The strategy is simple: we check whether additional parameters should be included in the model. Instead of the scatterplot matrix of the original data we look at the scatterplots of the (formal) residuals from this simple model. These scatterplots are called **added variable plots**.

To point out the difference to the scatterplot matrix of the original data: linear structures in the scatterplot the original data are a clear hint for linear dependency. Non linear structures, such as a triangular shape in some of the scatterplots can a reflect a corresponding dependency structure. But it can also be an artefact, based on the distribution and correlation structure of the regressors. In general, they do not have a simple interpretation. In contrast, the representations in the matrix of the added variable plots are adjusted for linear effects of the preceding variables. As a consequence, they depend on the order in which variables are taken into the model. They correct for linear effects resulting from the correlation of preceding variables. This avoids several artefacts, and the added variable plots are easier to interpret.

a:ch04:12

Exercise 4.6	
	<p>Modify the following function <code>pairslm()</code> so that it calculates the residuals of the regression of all original variables in matrix <code>x</code> by regression after the new variable <code>x\$fit</code> and produce a scatterplot matrix of these residuals.</p> <pre><code>pairslm &lt;- function(model, x, ... )   { x\$fit &lt;- lm(model, x)\$fit; pairs(x, ...)}</code></pre> <p>Add title, legends etc..</p> <p>Use the "trees" data set as an example.</p>

We discussed the transition from  $p'$  to  $p' + 1$  variable. The scatterplot matrix allows to get a quick overview of a (not too large) number of candidates (three potential additional regressors in our case). The transition from  $p$  to  $p - 1$  for the elimination of a variable, is in a certain sense dual to this. It corresponds to stepwise elimination, the other basic strategy.

Instead of using individual variables step by step, it can be more efficient to select linear combinations of several variables and use these synthetic variables in a model. Methods to achieve this are discussed under the heading **principal component analysis** and supported by function `prcomp()` in `library(mva)`. We will return to this in a later example (page 4-51).

The experience with the linear models teaches us, that the marginal relations are only half of the truth. Instead of using the regressors individually, we have to use a stepwise orthogonalisation. Component-wise interpretation of estimators becomes questionable. The estimation may depend on the order in which regressors are used.

In more complex situations, formal methods may lead astray. Technical skills are needed. Unfortunately knowledge about how interventions affect the validity of formal methods is rather limited. This makes it more important to judge chosen strategies critically using simulations.

#### 4.7.2 Non-linear Case

Non-linear relations in higher dimension pose a challenge. Besides a good collection of methods we need a repertoire of examples showing which structures can occur and where we have to pay attention. The following example, the cusp singularity is one of these examples: it is one of the most simple structures that can occur in higher dimensions. The base is here a two dimension structure, a surface, that is not trivially imbedded in a three dimension space. The interesting feature here is the bifurcation from a unimodal to a bimodal situation.

##### *Example: Cusp Non-linearity*

The most simple example can be illustrated with reference to physical applications. In physical systems, probability distributions often are related to energy levels; (local) Minima of the energy correspond to modes of the distribution. A typical relation is: if the energy behaves as  $\varphi(y)$ , after standardisation the distribution behaves like  $e^{-\varphi(y)}$ . If  $\varphi(y)$  is quadratic at the minimum, up to scale transformation, we get distributions from the normal distributions family.

Differential topology teaches us that the qualitative image persist under small perturbations or variations. At least locally the energy stays approximatively quadratic, and the normal distributions stay at least approximatively a suitable distribution family.

The behaviour changes drastically if the behaviour behaves locally like  $y^4$ . Small perturbations can make the potential locally a quadratic function. But they can as well lead to breaking the local minimum into two minima. The typical image has the form

$$\varphi(y; u, v) = y^4 + u \cdot y^2 + v \cdot y. \quad (4.3)$$

f:cusp0

The variations are controlled by the parameters  $u, v$ . A dynamic interpretation helps to understand the situation: imagine that  $u, v$  are external parameters which can change. We know this scenario from magnetic hysteresis:  $y$  gives the magnetisation in a certain direction,  $u$  takes the role of temperature;  $v$  that of an external magnetic field. At high temperature, magnetisation follows the external magnetic field directly. At lower temperatures, the material shows memory. The magnetisation does not only depend on the outer magnetic field, but also on the previous magnetisation.

We know similar “memory effect” in other areas. Imagine a market with prices  $y$ , costs  $v$  and a “pressure of competition”  $u$ . With sufficient competition the prices will follow the cost, more or less, given other conditions unchanged. In a monopoly situation, prices seem to have a memory: once they have risen, they will only become lower if the costs reduce drastically.

The “unfolding” of the potential  $y^4$  given in formula 4.3 has a characteristic form. From

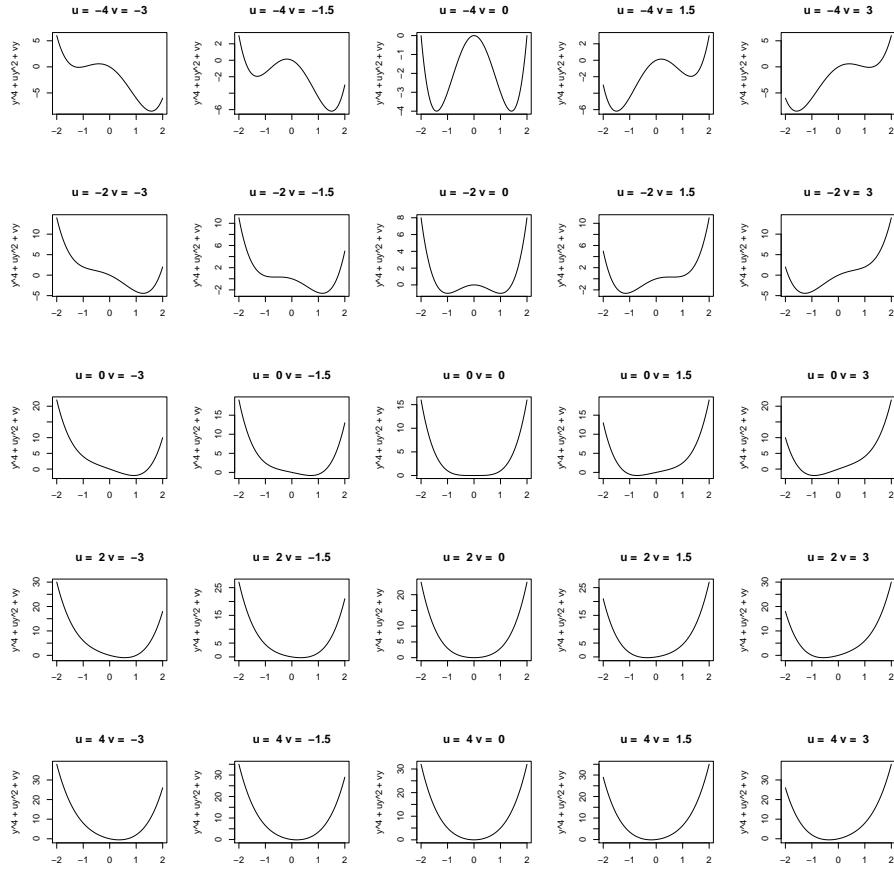
$$\varphi'(y; u, v) = 4y^3 + 2u \cdot y + v = 0 \quad (4.4)$$

f:cusp1

we get the critical points (see Abb. 4.3). See figure 4.3 on page 4-33.

fig:ch04cuspsurf1

4-32

DIMENSIONS 1, 2, 3, ...,  $\infty$ Figure 4.2 Unfolding of  $y^4$ :  $\varphi(y; u, v) = y^4 + u \cdot y^2 + v \cdot y$ 

A projection to the  $u, v$  gives a cusp (Fig. 4.4). For parameter values in the interior of this cusp the potential has two local minima, outside of the cusp there is only one extreme value.

**ToDo:** adjust scale, orientation

**ToDo:** 2/3

ch04cuspine

The distributions corresponding to these potentials are – up to scale transformation for normalisation –

$$p(y; u, v) \propto e^{-(y^4 + u \cdot y^2 + v \cdot y)}.$$

(4.5) propy4exp

The structure of the potentials carries over to the corresponding distributions.

The situation appears mostly harmless: the parameter space (the space of the regressors)  $x = (u, v)$  has only two dimensions. The distribution is one dimensional with a smooth density. But

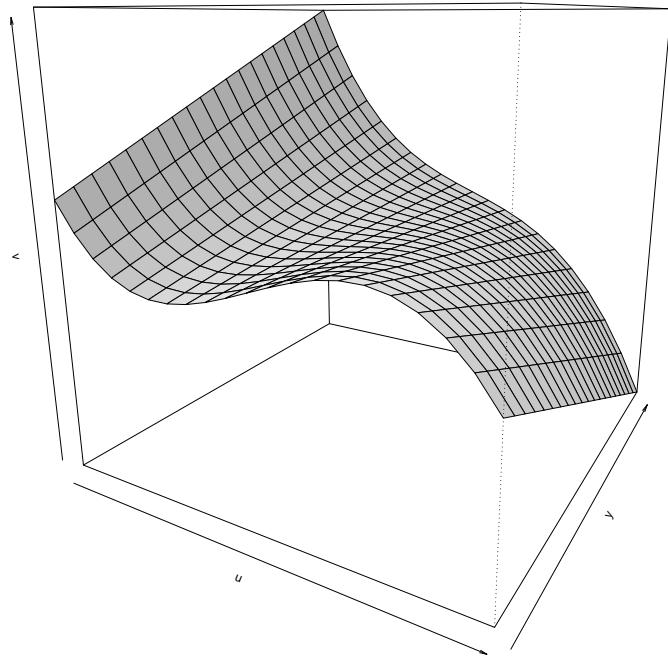


Figure 4.3 Critical points  $\varphi'(y; u, v) = 4y^3 + 2u \cdot y + v = 0$

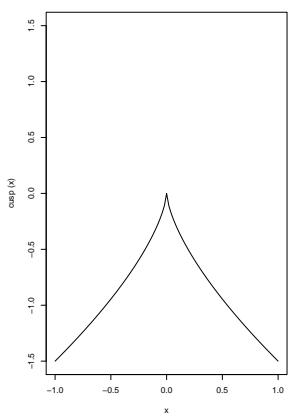
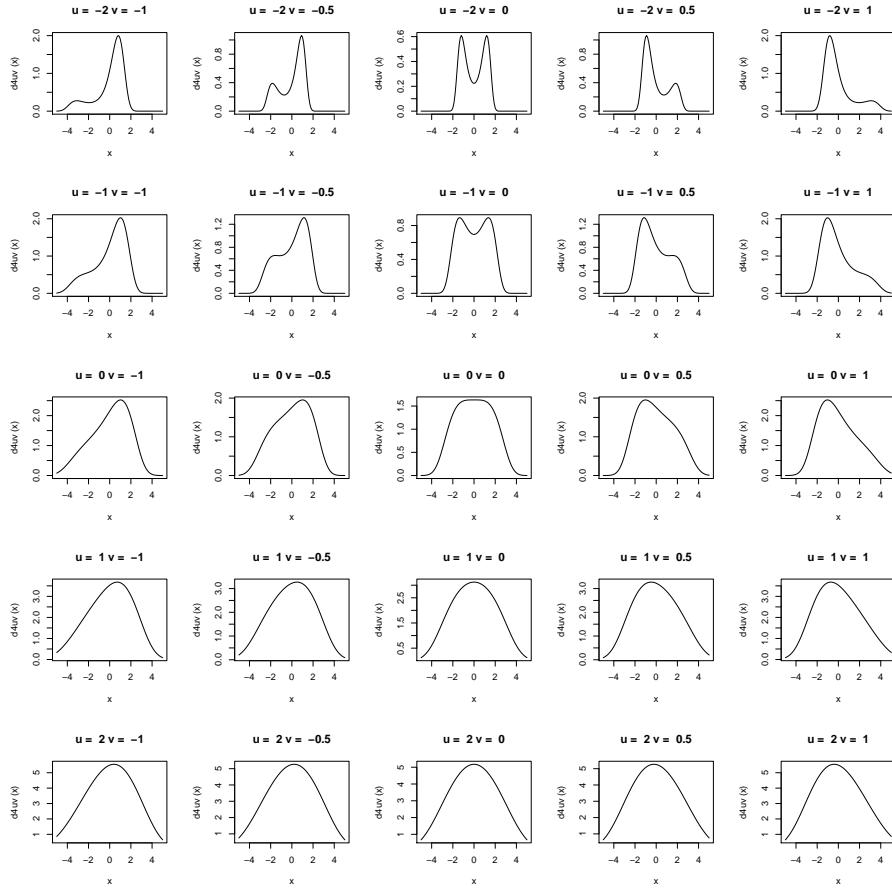


Figure 4.4 Boundary between uni- and bimodality in  $(u, v)$  space

4-34

DIMENSIONS 1, 2, 3, ...,  $\infty$ Figure 4.5  $p(y; u, v) \propto e^{-(y^4 + u \cdot y^2 + v \cdot y)}$ 

the situation cannot be captured sufficiently by linear methods. The typical non-linear effect will not be recognised if you are not prepared for it. Only an overall picture in three dimensions will convey the underlying structure.

a:ch04:13

This simple example is a challenge. How can we diagnose a structure of this kind?

Exercise 4.7	
	Write a function <code>dx4exp(x, u, v)</code> to calculate the centred probability densities for (4.5). You have to integrate the density from (4.5) to determine the norming constant and to calculate the expected value to centre the density. For both, you can use a numerical integration with <code>integrate()</code> .
	(cont.)→

Exercise 4.7	(cont.)
***	<p>For values <math>u, v</math> on a grid in <math>u = -2 \dots 2</math> and <math>v = -1 \dots 1</math> simulate 100 random numbers from <math>\text{dx4exp}(x, u, v)</math>. Use the methods from chapter 2 to analyse the data.</p> <p>Can you detect hints indicating the not-linear structure?</p> <p>Is the bimodality recognisable?</p> <p>How much of the structure can you identify?</p>

In non-linear relations, joint dependency can have a large importance. In general this requires prudence in modelling. Non-linear relations can be hidden in projections. Artefact of the (linear) projection can give an image that does not correspond to the original relations.

#### 4.7.3 Melbourne Temperature Data

R. Hyndman pointed out the bifurcation to bimodality in the Melbourne temperature data set [Hyndman1996Estimating-and-](#)<sup>[T2]</sup>. We use an extended version of the data set<sup>2</sup> and analyse the day by day difference in temperature at 15h (the daily report reference time) conditioned on today's temperature and pressure at the reference time.

We are interested in modality, and classical methods like those introduced in chapter 1 are not tuned for uncovering multimodality. One reason is that <sup>ch:01</sup>modality is a multi scale problem. If we use for example kernel density estimation (see section 1.3), typically we get exactly one mode for sufficiently large bandwidth, and a mode at every data point if the bandwidth is sufficiently low. To understand modality we have to look at different scales, and find the suitable once.

Instead of going into this problem of bandwidth choice, we use the shorth plot, which allows a multi scale representation. For each data point  $x$ , the plot shows the minimum length of an interval containing  $x$  and covering at least a proportion  $\alpha$  of the data. Using this plot, we can investigate several levels of  $\alpha$  simultaneously. For details and source code, see <http://lshorth.rforge.r-project.org/>.

The shorth plot view is in figure 4.6. The full picture reveals a cusp-type bifurcation. Figure 4.7 shows the shorth plot for the temperature difference at 15h (the daily report reference time) to next day's temperature, conditioned on today's temperature and pressure. It reveals the modality split as well as the skewness which is pressure dependent.

**ToDo:**  
Melbourne  
data pub-  
lic?

<sup>2</sup> Melbourne temperature data 1955-2007, provided by the Bureau of Meteorology, Victorian Climate Services Centre, Melbourne.

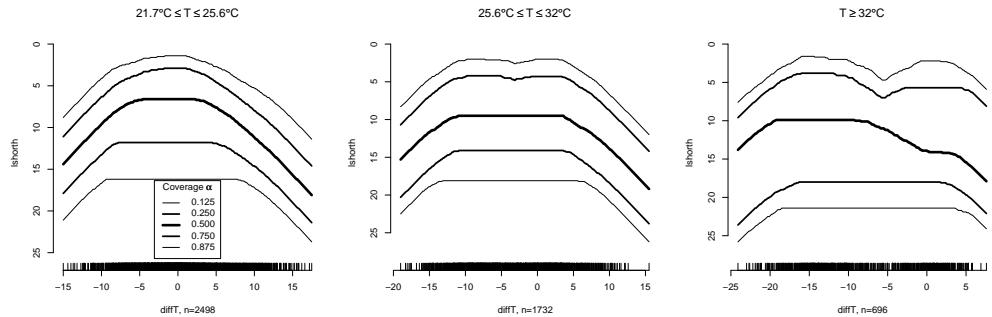


Figure 4.6 *Shorth plot at coverage levels  $\alpha = 0.125, 0.25, 0.5, 0.75, 0.875$  for Melbourne day by day temperature difference at 15:00h conditioned at today's temperature. A bifurcation to bimodality occurs at high temperatures.*

fig:melbournet

## HIGHER DIMENSIONS

4-37

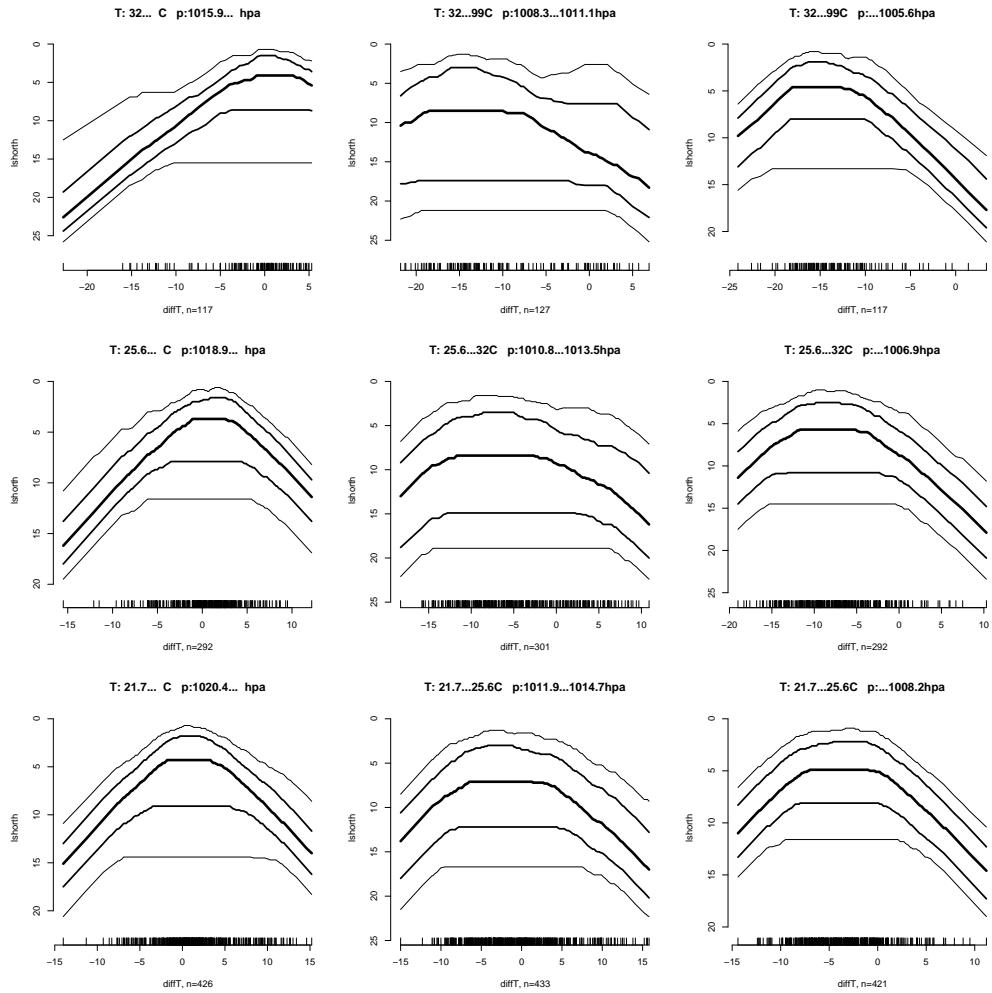


Figure 4.7 Plot matrix of shorth plots for varying temperature ranges (vertical) and varying pressures (horizontal). Shorth plots at coverage levels  $\alpha = 0.125, 0.25, 0.5, 0.75, 0.875$  for Melbourne day by day temperature difference at 15:00h, conditioned at today's temperature and pressure.

fig:melbourne

#### 4.7.4 “Curse of Dimensionality”

With the lack of adapted co-ordinate systems, an extensive search for interesting projections or sections is necessary. The number of possibilities increases rapidly with the dimension. For illustration: to identify a cube, at least the corner points must be detected. In  $d$  dimensions there are  $2^d$  corner points. The number increases exponentially with the dimension. This is an aspect of the problem known as ***curse of dimensionality***.

Seen from another point of view: if we look at the data points which are extremal in at least one variable dimension, in the one dimensional case we get two points. In  $d$  dimensions typically we get  $2^d$  extremal points. If we do not look at co-ordinate directions, but in arbitrary directions, typically any point is extremal, if  $d$  gets large.

A third aspect: in  $d$  dimensional space almost every point is isolated. Localisations, like those discussed in section 2.5, will break down. If we take a neighbourhood around some point covering a proportion  $p$  of the data range, for example  $p = 10\%$  of the span of the variables, in one dimension typically we cover a proportion of  $p$  of the data points. In  $d$  dimensions this is only a fraction in the order of  $p^d$ . So for example in 6 dimensions we need several million data points in order to have non-empty environments for most data points.

#### 4.7.5 Case Study: Crabs

##### **To Do:**

give  
proper  
comment  
to external data  
access  
**To Do:** adjust style

The Potamon crab data set <http://statlab.uni-heidelberg.de/data/khabur/khabur.data> contains morphological data of 234 specimens from the collection of the Senckenberg museum.

In each specimen, a set of morphometric variables was measured: length and width of the body (the carapace), of the scissors (chelipeds) and one of the walking legs (pereiopods), resulting in 21 basic measurements (see figure 4.8). As additional information were recorded: the sex (male, female, juvenile) and the location of origin.

The specimens had been collected in various regions: Jordan, Orontes, Euphrates, Khabur (a tributary of Euphrates), Italy and Greece. The specimens from Italy (*Potamon fluviatilis*) and the northern Greek mainland (*Potamon ibericum*), were included as reference material, clearly characterised by their gonopods.

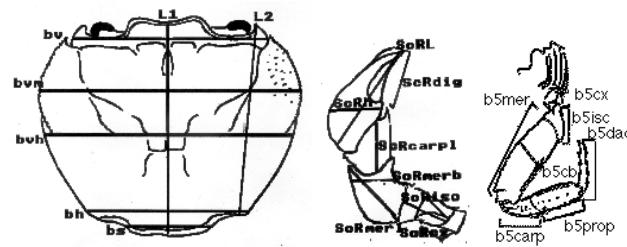


Figure 4.8 *Crab morphological data*

**fig:crabmorph**

---

**Input**  
`khabur <- read.delim(file="..../data/khabur.data")`

*Strategy*

It would be tempting to analyse the material using sophisticated methodology, like deformable templates [1]. However, adequate software was not existant, and even if it would have been, convincing arguments to communicate any findings would not have been easy. Thus, we headed for a more down-to-earth approach. As a guide, we used a (multivariate) linear model  $Y_{id,feature,repetition}$ . This does not mean that we apply a linear model immediately to our data without their preparation, but we use the linear model as a guiding example. In data analysis, we should make use of the experience built up in this classical area. Linear models are well understood in statistics, and we could use them as a guide to avoid the pitfalls and errors which have been well studied in the past.

Part of statistical analysis is the communication of results. The communication process can be inhibited by the use of obscure methods. We headed for simple “classical” methods, possibly loosing some power, but maintaining the chance to check the assumptions and to interpret the results critically.

With 21 basic parameters for 234 specimens the material is so rich that almost any approach will find a significant result, but possibly of dubious relevance.

A classical way to fatal errors would be to build a hypothesis on the same information base that is used to confirm it. The nominal significance level of any test assumes given model assumptions and hypothesis. If we derive a hypothesis on the basis of the data material we use for testing, the hypothesis (and model) itself is a stochastic entity. The nominal significance level needs to be corrected for possible dependence. In general, this is not feasible. The classical way out, well known, but rarely used, is to separate the data material into a training set and an evaluation set, using a separation strategy which enforces stochastic independence of training set and an evaluation.

Formally, we would be free to use the training set to build our model and hypothesis as long as we use a stochastically independent subsample for evaluation. To avoid running in circles, we adopted a strict strategy during the modelling process as well. We discarded any information we had used in one decision step if it might have affected subsequent decisions. As a jargon, we said that information is used up. This may be over-conservative, but it puts us on the safe side.

Another prominent pitfall would be an implicit inclusion of a separation into groups. The background problem can be seen as a classification problem. Following this track using a clustering or discriminant approach would “confirm” the existence of classes, but would not allow to check the underlying assumption. So the results would be worthless. We took great care not to assume any separation into a priori classes.

*Data Definition*

In all statistical analyses, in particular in biometrical applications, a proper definition of the variables is necessary. This implies a proper definition of variable scales and categories. For example, “left” and “right” are nominal categories, but in a biometrical context they may be associated with a role as dominant or subdominant side. If the chirality is not fixed, we may have a subpopulation where “left” and “right” would need to be exchanged. To make things more complicated, crabs, like other decapods, can regenerate limbs. So, an animal with a larger right cheliped may loose the limb at some time, regenerating it later and keeping the left as the dominant one. For those variables which have been taken on the right and left side, extra checks

are needed for proper handling of chirality. For the morphometric variables, the definitions are based on pragmatic considerations. The morphometric variables are simplified means to describe a complex 3 dimensional structure. To qualify as a variable a feature must have a clear morphometric definition, and it must be accessible for measurement in a reasonable way. In growth processes and morphometric studies the proper scale for an analysis may be defined by a "growth curve" and if these growth curves are known, proper adjustment may be taken. For the Freshwater crabs of the genus *Potamon*, we could not rely on literature information on the growth pattern, so appropriate transformations had to be derived from the data material itself.

### *Data Preparation*

As usual, the first step was checking and cleaning the data. All data had been registered using an electronic caliper which feeds the data online into a data base, without any intervention. By common laboratory praxis, the data set was considered correct and error free. By common statistical praxis, and guided by the multivariate linear model as a leading example, we knew that outliers might critically affect the analysis and made a routine check for outliers. The analysis pointed at several unpleasible values.

A more thorough outlier check was added. Although we had only 3 repeated measurements per item, the sample size of 234 specimens allows an estimation of the variance of the measurement error,  $Var(Y_{.,feature,..})$ , assuming a common variance per feature using the model  $Y_{id,feature,repetititon.} = \mu_{id,feature} + \varepsilon_{id,feature,repetititon.}$  Assuming a  $\chi^2$ -distribution, a warning limit for the individual measurement error variance  $Var(Y_{id,feature,..})$  can be derived and used to select those registrations needing a closer inspection.

We chose a warning limit for an error rate of 5%. A total of 15% of the measurements needed a closer inspection, following this strategy - much more than expected. The source of the problem was identified in the caliper used for the morphometric values. If the registration button was pressed too long, the caliper would register a value repeatedly, overwriting cells which were assigned to subsequent values - an effect which had passed unobserved before.

We considered the information of the repetitions used up by this cleaning step and used the mean over the repetitions only later on.

### *Training Set*

The next problem we addressed was splitting the data material into an exploratory, or training set and a confirmation, or evaluation set. The decision to take was, which proportion of the data to allocate to each. When separating a data set into a training set and an evaluation set, the size decision is in a dilemma. For modelling, you want as many data as possible to be able to detect a structure. More data in the training set may lead to more specific models, hence possibly more powerful tests in the evaluation step. But any data put into the training set reduces the sample size of the evaluation set, hence reduces the power of any later test. Reserving data for the training set can only be justified if the loss of power on the one hand is compensated by better modelling and a stricter definition of hypothesis on the other side.

Lacking a proper statistical criterion to determine the size of the training set, we took an experimental approach. We intended to use the training set for interactive graphical data analysis to find appropriate models and hypotheses for a later evaluation. To determine an appropriate sample size, we used the variance/covariance structure from the complete data set. We generated a series of pseudo-random samples from various distributions and related by various power

laws, matching the variance/covariance structure of our data set. For each pseudo random sample we presented a scatter plot matrix which we inspected visually to find out a sample size which was big enough to make the most common models apparent. This series of experiments suggested a sample size about 20 to be the minimum useful size.

Keeping in mind linear models we knew that if we would make any ANOVA type analysis later, a balanced design would make analysis much easier. The number of specimen per location varied from 40 to 50 for the locations of interest. We made the arbitrary decision to reserve 30 animals for the evaluation data set, and use the additional ones for the training data set, using a randomized selection.

### *Preliminary Analysis*

For a preliminary analysis, only the training data part of the data set was used. A first exploratory analysis on the data in their original form showed no indication of any reasonable grouping. There was a clear difference in overall size between regionally separated groups, seen most clearly in Box & Whisker plots. But apart from this, scatter plot matrices or other tools did not reveal any particular information apart from a strong correlation of obviously related variables (see fig. 4). In particular, no difference between the various regional groups was visible in other respects than overall size.

The next step was to apply classical methods as exploratory tools. To avoid running into too many dimensions, we separated the variables into groups relating to the carapace, the chelae and to the legs. For each of these variable groups, we applied a principal component analysis to the training data set. At this point, we had no information whether the underlying assumptions of principal component analysis on the joint distribution would hold, and had no information about the reliability of the result. We used PCA as a speculative exploratory tool only to get some hints for further modelling. (insert table 1) Table 1 shows the results of a principal component analysis, applied to the carapace traits. A formal application of a principal component analysis would identify the linear combination

$$F1 = -0.380 * L1 + (-0.379) * L2 \dots + (-0.371) * Bs$$

as dominating variance component. Instead of following a formal analysis, we did prefer to head for a model which allows a sensible interpretation. The formal first principal component puts nearly equal weight to all variables. This component can be interpreted as the general size of the animal. Instead of using the formal linear combination, we can pick any of the variables without much loss in precision.

At this point we went back to the subject information at hand and asked which of the variables would give the most reliable measurement of the size. By morphology and by the set up of the measurement process, the overall length  $L1$  has the most precise definition and can be measured most reliably. Instead of using the formal combination  $F1$ , we made the decision to use the overall length  $L1$ . All other carapace variables were standardized by  $L1$ . Following our general strategy we considered the information of  $L1$  used up after this standardization and unusable for any further analysis.

An additional principal component analysis on the length standardized data differs only slightly from the results given in the table above. The next principal component is marked by the contrast between the caudal carapacial width  $Bh$  and the width at the pleon base  $Bs$  on the one hand and the rest of the carapace traits. So it reflects the relative width of the central part of the body. Next, the principal component analysis draws our attention to the epibranchial

and pleon base width  $Bv$  and  $Bs$ , contrasted with the other measurements. So it catches an assumption of how strongly the body is pointed at the extremes. The next components reflect the general shape ratio (length/width) and the relative width at the middle carapace  $Bvm$ ,  $Bvh$  contrasted with the other variables, respectively.

Like for the carapace, the principal component analysis for the cheliped and the leg led to a morphologically sensible rescaling and grouping of the variables.

As a side effect, the standardization suggested by principal component analysis leads to a separation of the data material with data from different regions of origin isolated as clearly distinct groups. This is already visible in a bivariate plot, illustrated in figure 5 which shows two variables for the complete data set. It should be pointed out that the information about groups by origin has not been used in the definition of the rescaling. The grouping apparent in figure 5 appears natural and has not been induced by an a priori classification with reference to geographic origin.

Only information from the training set has been used to define the scaling and to find the hypothesis that we have distinct groups. The main part of the data material is unused and is available for an independent validation.

### *Final Analysis*

Using the training data set, several hypothesis concerning differences between groups could be formulated. The details are published elsewhere. In principle, the remainder of the analysis is classical statistics: the usual assumptions about joint distributions can be inspected using the training set. Nothing suspicious was found for the standardized data. Once the hypothesis has been formulated using an independent source of information, the training data set, the hypothesis can be checked on the evaluation set.

Finer points need to be taken into account. In the preliminary analysis, we separated information about carapace, chelipeds and legs. Morphological variables related to these groups of characters will show stochastic dependence, and this dependence is prone to affect any analysis. To correct for this dependence, we used a Bonferroni adjustment. Another stochastic dependence arises from generalised hypotheses. On the basis of the training set we could specify the variables which might show a difference between groups, but we could not specify which groups would have significant differences. For a formal test, we would like to identify the differing groups, once we have established that certain variables show a difference at all. We used Scheffé's post hoc test to give the proper corrections. With these precautions, the rest is classical multivariate analysis of covariance. A sample result is shown in tables 2 and 3.

**ToDo:**  
add tables

### *Discussion*

Statistical analysis must end with the formal evaluation. Several differences between the groups have been established as statistically significant. In particular, next to the groups from Greece and Italy, the Khabur group, which initially had been considered as merged with the Jordan/Euphrat/Orontes groups, turned out as a well separated group with a distinct morphology in all examined characters. The other groups Jordan, Orontes and Euphrates show high similarities to each other as they can either only be separated in single characters like chelae or walking legs or they can not be separated at all. However as mentioned above the Euphrates group will not be discussed here as the group is not homogenous including several species. This material can not be treated as a single group and results are not comparable to the other analysed groups.

Additionally, all data sets show the extremely high variability and heterogeneity of this material and thus, a new detailed analyse especially of the groups from Iraq and Iran is necessary to get meaningful results from these groups. At this point, biology takes over again. At first a comparison of the statistical results with the morphological data is necessary. Former morphological results are primarily based on characters of the copulatory system, which are relevant for a successful reproduction.<sup>[4]</sup> These characters can be used for a delimitation of species, as it has been done by<sup>[4]</sup>. Generally, the groups separated by statistical methods correspond perfectly to the morphological groupings proposed in the aforementioned revisionary papers. The subgenera can be perfectly delimited by morphometric methods. Within the subgenera only in Potamon (Potamon) two species were compared to each other. These are very similar and can only be distinguished by gonopodial characters and by morphometric differences of the 5th leg. This study shows full agreement of morphometric groupings and hierarchies with the ones already defined by the morphology of the reproductive system. Together with the congruence of the geographic distribution of the taxa with palaeogeographic settings it is evident that the system of freshwater crabs of Europe and the Middle East proposed by<sup>[4]</sup> is the most congruent in all aspects so far examined. This study also shows that no single qualitative difference of carapace and limbs can be defined as a diagnostic character for taxa. The overlap is such that only morphometric methods with advanced statistics can help to detect distinguishable groupings. It makes, therefore, no sense to use such characters with a broad overlap as basis for taxonomic groupings or attribute any phylogenetic significance to them. We confirm thus, the usefulness of a tradition using the morphology of the reproductive system as a prime character in the systematics of potamid freshwater crabs. The morphometry may prove useful as a supplementary tool, if it can be confirmed by further studies on a much broader scale that it really results in meaningful and stable groupings as suggested by the present case study. Thus, however is still to prove.

#### 4.7.6 Case Study: Body Fat

As a running example we use here the fat data set. This data set is published repeatedly in the literature and in R among others available in library *UsingR*.

Target of the investigation related to this data set is the determination of the proportion of body fat. The most reliable method is to use a water bed to determine the average density of the tissue and to infer from the density to the body fat proportion. This measurement is rather intricate. Can it be replaced by body parameters that allow a more simple measurement? The parameters at disposition are summarised in table 4.16.

The survey in table 4.16 tells us that metric variables and US scales are mixed. To make interpretation easier for us, we convert all information to metric units.

---

*Input*

```
library("UsingR")
data(fat)
fat$weightkg <- fat$weight*0.453
fat$heightcm <- fat$height * 2.54
fat$ffweightkg <- fat$ffweight*0.453
```

The variables *body.fat* and *body.fat.siri* are derived from the observed values of *density*. These formulas reflect assumptions about the mean density of fat and of fat free tissue. With these assumptions the fat proportion can be calculated (or rather: estimated) from the *density*.

**ToDo:**  
reduce  
body fat  
case study  
drastically

**ToDo:**  
reduce  
preliminary  
information on  
fat

Name	Variable	Unit, Remarks
case	Case Number	
body.fat	Percent body fat using Brozek's equation, $457/Density - 414.2$	
body.fat.siri	Percent body fat using Siri's equation, $495/Density - 450$	
density	Density	$[g/cm^2]$
age	Age	$[yrs]$
weight	Weight	$[lbs]$
height	Height	$[inches]$
BMI	Adiposity index = $Weight/Height^2$	$[kg/m^2]$
ffweight	Fat Free Weight = $(1 - fraction of body fat) * Weight$ , using Brozek's formula	$[lbs]$
neck	Neck circumference	$[cm]$
chest	Chest circumference	$[cm]$
abdomen	Abdomen circumference "at the umbilicus and level with the iliac crest"	$[cm]$
hip	Hip circumference	$[cm]$
thigh	Thigh circumference	$[cm]$
knee	Knee circumference	$[cm]$
ankle	Ankle circumference	$[cm]$
bicep	Extended biceps circumference	$[cm]$
forearm	Forearm circumference	$[cm]$
wrist	Wrist circumference "distal to the styloid processes"	$[cm]$

Table 4.16 Fat data set: variables

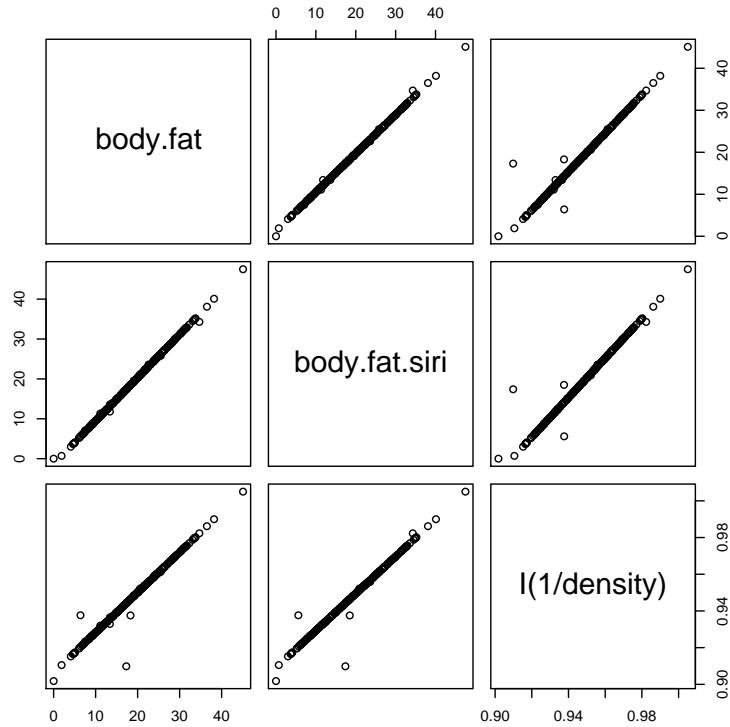
**fatvars**

Both formulas use a density dependent factor  $1/density$ . Up to (given or assumed) constants is this the relevant term for us (and not *density*).

The first step is a critical inspection and cleaning of the data set. This is almost always necessary, not only for higher dimensional data sets. In higher dimensional data sets though we often have redundancies that allow for consistency checks and possibly for corrections. In our case, the variables *body.fat*, *body.fat.siri*, *ffweight* and *BMI* are derived variables that are in a deterministic relation to other variables.

We first look at the group *body.fat*, *body.fat.siri*,  $1/density$ . The pair-wise scatterplots should show a straight line. *pairs()* provides good service. We use the formula based notation here. To mark that  $1/density$  should be calculate, and the / operator is not to be understood as a formula operator, we have to mark up this term.

*Input*  
`pairs(~body.fat + body.fat.siri + I(1/density), data = fat)`



The inconsistent values and outliers are obvious. Unfortunately in R it is not simple to mark values in a scatterplot matrix.

a:ch04:14

Exercise 4.8	
	<p>Use functions <code>plot()</code> and <code>identify()</code>, to generate the following output:</p> <div style="display: flex; justify-content: space-around;"> <div style="text-align: center;"> <p>body fat</p> <p>1/density</p> </div> <div style="text-align: center;"> <p>body fat siri</p> <p>1/density</p> </div> </div> <div style="display: flex; justify-content: space-around; margin-top: 20px;"> <div style="text-align: center;"> <p>ffweightkg</p> <p>(1-body.fat/100)*weightkg</p> </div> <div style="text-align: center;"> <p>BMI</p> <p>weightkg/(heightcm/100)^2</p> </div> </div>

If an obvious correction is possible, it should be done here. Of course the correction needs to be noted in the report of the analysis. case 42 is simple: a height of 0.73m for a weight 63.5kg is not plausible and not consistent with BMI 29.9. From the BMI we can recalculate the height. Instead of the value of 29.5 inch the height should presumably be 69.5 inch.

*Input*

```
fat$height [42] <- 69.5
fat$heightcm[42] <- fat$height[42] * 2.54
```

Case 216 is discretionary. The density is extremely low, the BMI extremely high. On the other hand, the body measurements fit to these extremes. This case can be an outlier which would distort the analysis. But it could also be an observation which is particularly informative. We note it as a particularity.

After this preliminary inspection we clean up the data set. We remove the variables that contain no additional information or have been replaced. As a target variable, we use `body.fat`. However we keep the variable `density` for later purposes.

---

*Input*

```
fat$weight <- NULL
fat$height <- NULL
fat$ffweight <- NULL
fat$ffweightkg <- NULL
fat$body.fat.siri <- NULL
```

Some indices are in common use to describe the body constitution (see Figure 4.9). Some time ago, the rule of thumb was ‘optimal weight [kg] = height [cm] -100’. Today the “body mass index”  $BMI = \text{weight}/\text{height}^2$  is usual. (Commercially available scales determine the body fat by electrical impedance. This variable is not contained in the fat data.)

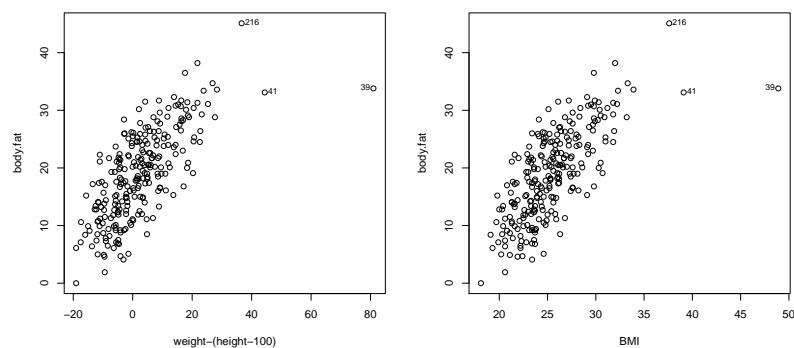


Figure 4.9 *Body fat against conventional indices*

**fat01scat**

We can use the conventional indices in a linear model. We exclude the obvious outliers and possible leverage points. This can be achieved with the `subset` parameter of `lm()`.

For the rule of thumb ‘optimal weight = height -100’ we get:

---

*Input*

```
lm.height <- lm(body.fat~I(weightkg-(heightcm-100)),
  data = fat,
  subset = -c(39, 41, 216))
summary(lm.height)
```

---

*Output*

Call:

```
lm(formula = body.fat ~ I(weightkg - (heightcm - 100)), data = fat,
  subset = -c(39, 41, 216))
```

Residuals:

Min	1Q	Median	3Q	Max
-11.90734	-3.68697	-0.05303	3.65458	12.28000

Coefficients:

```

Estimate Std. Error t value
(Intercept) 17.70722 0.33296 53.18
I(weightkg - (heightcm - 100)) 0.54557 0.03283 16.62
Pr(>|t|)
(Intercept) <2e-16 ***
I(weightkg - (heightcm - 100)) <2e-16 ***
---
Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

```

Residual standard error: 5.166 on 247 degrees of freedom  
 Multiple R-squared: 0.5279, Adjusted R-squared: 0.526  
 F-statistic: 276.2 on 1 and 247 DF, p-value: < 2.2e-16

The regression of *body.fat* for *BMI* results in:

---

<pre> lm.BMI &lt;- lm(body.fat ~ BMI,                data = fat,                subset = -c(39, 41, 216)) summary(lm.BMI) </pre>	<i>Input</i>
--	--------------

---

<pre> Call: lm(formula = body.fat ~ BMI, data = fat, subset = -c(39, 41, 216))  Residuals:     Min      1Q  Median      3Q     Max  -12.49460 -3.53561 -0.05228  3.69129 11.72720   Coefficients: Estimate Std. Error t value Pr(&gt; t )     (Intercept) -25.6130    2.6212 -9.772 &lt;2e-16 *** BMI         1.7564    0.1031 17.042 &lt;2e-16 *** --- Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1 </pre>	<i>Output</i>
--	---------------

Residual standard error: 5.097 on 247 degrees of freedom  
 Multiple R-squared: 0.5404, Adjusted R-squared: 0.5385  
 F-statistic: 290.4 on 1 and 247 DF, p-value: < 2.2e-16

The fit of  $R^2 = 0.53$  resp.  $R^2 = 0.54$  is poor in both cases.

Even using all data points and all regressors, we get at most  $R^2 = 0.75$ :

---

<pre> lm.fullres &lt;- lm(body.fat ~ age + BMI + neck + chest + abdomen + hip + thigh + knee + ankle + bicep + forearm + wrist + weightkg + heightcm, data = fat) summary(lm.fullres) </pre>	<i>Input</i>
--	--------------

```
Call: lm(formula = body.fat ~ age + BMI + neck + chest + abdomen +
       hip + thigh + knee + ankle + bicep + forearm + wrist + weightkg +
       heightcm, data = fat)

Residuals:
    Min      1Q  Median      3Q     Max 
-10.0761 -2.6118 -0.1055  2.8993  9.2691 

Coefficients:
            Estimate Std. Error t value Pr(>|t|)    
(Intercept) -50.804727  36.489198 -1.392  0.16513  
age          0.061005  0.029862  2.043  0.04217 *  
BMI          0.782993  0.733562  1.067  0.28688  
neck         -0.439082  0.218157 -2.013  0.04528 *  
chest        -0.040915  0.098266 -0.416  0.67751  
abdomen       0.866361  0.085550 10.127 < 2e-16 *** 
hip           -0.206231  0.136298 -1.513  0.13159  
thigh         0.246127  0.135373  1.818  0.07031 .  
knee          -0.005706  0.229564 -0.025  0.98019  
ankle         0.135779  0.208314  0.652  0.51516  
bicep         0.149100  0.159807  0.933  0.35177  
forearm        0.409032  0.186022  2.199  0.02886 *  
wrist          -1.514111  0.493759 -3.066  0.00242 ** 
weightkg      -0.389753  0.221592 -1.759  0.07989 .  
heightcm       0.187196  0.199854  0.937  0.34989 

---
Signif. codes:  0 ‘***’ 0.001 ‘**’ 0.01 ‘*’ 0.05 ‘.’ 0.1 ‘ ’ 1

Residual standard error: 3.991 on 237 degrees of freedom
Multiple R-squared:  0.7497,    Adjusted R-squared:  0.7349 
F-statistic: 50.7 on 14 and 237 DF,  p-value: < 2.2e-16
```

This is a model with 15 coefficients. The model is so complex that it is hard to interpret, and one would try to reduce the model. Instead of looking “by hand” for simpler models, this process can be automated. This is done for example by function `regsubsets()` in `library(leaps)`. The quadratic error (resp. the coefficient of determination  $R^2$ ) has to be modified: the quadratic error is trivially minimised if we take all regressors into the model, so the complete model will always be the optimal model in terms of the coefficient of determination. For model selection one uses variants of the quadratic error (resp. the coefficient of determination  $R^2$ ) that are adjusted for the number of parameters.

a:ch04:141

Exercise 4.9	
*	<pre>Use library(leaps) lm.reg &lt;- regsubsets(body.fat    age + BMI + neck + chest + abdomen + hip + thigh + knee + ankle + bicep + forearm + wrist + weightkg + heightcm, data = fat) and inspect the result with summary(lm.reg) plot(lm.reg, scale = "r2") plot(lm.reg, scale = "bic") plot(lm.reg, scale = "Cp") Hint: see help(plot.regsubsets)</pre>
*	Use the function <code>leaps()</code> for model selection.

But the tools which we discussed in chapter 2<sup>ch:02</sup> have tarnished now. The statistical statements in the summary are only valid, if model resp. hypothesis are stated independent of the data material. If a model is data based, the distribution of the estimated coefficients is unclear. We do not know how to define confidence intervals or how to perform tests. The software does not know that we are in a process of model selection and just returns the probabilities that apply for a fixed model under the assumption of a normal distribution.

The diagnostics, such as distribution plots of the residuals, becomes useless as well: the central limit theorem ensures that under weak independence conditions the residuals approximately have a normal distribution if we have a large number of terms, even if the normal distribution assumption does not hold for the errors themselves.

We are in a dead end.

We illustrate a different approach that goes further. For this we step back to the beginning of the analysis, after the preliminary inspection and data correction. To avoid running into the problem that the statistical distributions are influenced by preceding model selection steps, we split the data set. We use one part as a training set, used for model selection and to play with various alternatives. The rest is reserved as evaluation set. The information from this part is only used after model selection for the statistical analysis.

A closer look shows that the model selection is most critical for the estimation of the error, not so much for parameter estimation. If the error is estimated based on the data used for model selection, we tend to under-estimate the error. The evaluation part serves for a reliable estimation of the error and for residual diagnostics. This is a limited task. We reserve only a smaller proportion of the data for this.

---

*Input*

```
sel <- runif(dim(fat)[1])
fat$train <- sel < 2/3
rm(sel)
```

We discard the outliers from the training part.

---

*Input*

```
fat$train[c(39, 41, 216)] <- FALSE
summary(fat$train)
```

	Output	
Mode	FALSE	TRUE
logical	89	163

Our target variable is `body.fat`, or alternatively `1/density`.

We try to structure the variables with respect to the subject matter. For the density we have a definition from physics

$$\text{density} = \frac{\text{weight}}{\text{volume}}.$$

Among the variables that are candidate regressor we have just one variable that directly refers to weight (`weight` resp. `weightkg`), a series of variables, referring to body geometry, and `age`.

From the observed density and the observed weight we can retrieve the volume. We add this to our variable set. Since we have just one weight observation per person, there is no space for person related statistics.

---

*Input*

---

```
fat$vol <- fat$weightkg/fat$density
```

Next we try to model the volume `fat$vol` in terms of the more convenient variables on the body geometry. The body geometry reported are linear values. In a rough approximation, we can derive volume related variables. The only length information we have is in `height`. For the lack of better information we assume that that all part of the body have a length which is proportional to the height. We are aiming for a linear model. Hence linear factors can be neglected since they are estimated by the model anyhow. If we take cylindrical approximations for the parts of the body, up to linear factors we have

---

*Input*

---

```
fat$neckvol <- fat$neck^2 * fat$heightcm
fat$chestvol <- fat$chest^2 * fat$heightcm
fat$abdomenvol <- fat$abdomen^2 * fat$heightcm
fat$hipvol <- fat$hip^2 * fat$heightcm
fat$thighvol <- fat$thigh^2 * fat$heightcm
fat$kneevol <- fat$knee^2 * fat$heightcm
fat$anklevol <- fat$ankle^2 * fat$heightcm
fat$bicepvol <- fat$bicep^2 * fat$heightcm
fat$forearmvol <- fat$forearm^2 * fat$heightcm
fat$wristvol <- fat$wrist^2 * fat$heightcm
```

Next we inspect the internal structure of the volume oriented regressor candidates, using the training data set only. We use the function `prcomp()`, which for any given set of variables yields stepwise best linear predictors.

fatprincomp

For the approximative volumes of the body parts, the principal components are:

---

*Input*

---

```
pcfatvol <- prcomp(fat[, 20:29], subset = fat$train)
round(pcfatvol$rotation, 3)
```

	Output						
	PC1	PC2	PC3	PC4	PC5	PC6	PC7
neckvol	0.054	-0.017	0.063	-0.070	0.570	0.045	0.787
chestvol	0.548	0.457	0.691	0.012	-0.113	-0.004	-0.003
abdomenvol	0.650	0.290	-0.697	-0.059	0.049	-0.011	-0.017
hipvol	0.483	-0.748	0.099	0.437	-0.031	0.069	0.003
thighvol	0.185	-0.370	0.074	-0.870	-0.244	-0.017	0.079
kneevol	0.056	-0.082	0.055	-0.061	0.336	-0.822	-0.278
anklevol	0.016	-0.032	0.034	-0.005	0.140	-0.236	-0.047
bicepvol	0.051	-0.047	0.079	-0.180	0.552	0.508	-0.531
forearmvol	0.024	-0.019	0.074	-0.086	0.388	0.019	-0.104
wristvol	0.009	-0.005	0.017	0.002	0.110	-0.053	0.044
	PC8	PC9	PC10				
neckvol	-0.169	0.067	0.093				
chestvol	-0.012	0.005	0.002				
abdomenvol	0.016	-0.015	-0.002				
hipvol	0.018	0.011	0.004				
thighvol	0.012	-0.015	-0.019				
kneevol	-0.159	0.299	0.056				
anklevol	-0.120	-0.949	0.073				
bicepvol	-0.333	0.030	0.010				
forearmvol	0.906	-0.049	0.044				
wristvol	0.003	-0.048	-0.990				

The pattern of the signs for the loadings gives hints about the internal structure. The first principal component  $PC1$  is a linear combination of variables which essentially describes the torso. The second principal component contrasts the upper part of the torso (chest to abdomen) with the lower torso. The third contrasts the abdomen to the rest of the torso.

a:ch04:15

**Exercise 4.10** Draw a sketch that shows which body geometry features are represented by the next principal component  $PC_4, \dots, PC_{10}$ .

## ToDo:

discuss

scale dependency

ToDo: use

Q3: use principal components in modelling. pc regression?

```
lm.vol <- lm(vol ~ neckvol + chestvol + abdomenvol +  
    hipvol + thighvol + kneevol +  
    anklevol + bicepvol + forearmvol +  
    wristvol,  
    data = fat, subset = fat$train)  
summary(lm.vol)
```

1

```
Call:
lm(formula = vol ~ neckvol + chestvol + abdomenvol + hipvol +
    thighvol + kneevol + anklevol + bicervvol + forearmvol + wristvol)
```

```

data = fat, subset = fat$train)

Residuals:
    Min      1Q  Median      3Q     Max 
-5.8242 -1.1054  0.1403  1.1521  4.6173 

Coefficients:
            Estimate Std. Error t value Pr(>|t|)    
(Intercept) 2.723e+00  1.367e+00  1.991  0.04823 *  
neckvol     2.677e-05  9.456e-06  2.831  0.00527 ** 
chestvol    9.082e-06  1.260e-06  7.205  2.52e-11 *** 
abdomenvol  1.201e-05  1.129e-06 10.641 < 2e-16 *** 
hipvol      5.379e-06  1.766e-06  3.045  0.00274 ** 
thighvol    1.805e-05  3.284e-06  5.495  1.61e-07 *** 
kneevol     6.338e-06  9.379e-06  0.676  0.50023    
anklevol    2.579e-05  1.009e-05  2.556  0.01157 *  
bicepvol   2.360e-05  8.117e-06  2.908  0.00419 ** 
forearmvol  2.224e-05  9.179e-06  2.423  0.01655 *  
wristsvol  -1.775e-05  3.937e-05 -0.451  0.65280    
---
Signif. codes:  0 ‘***’ 0.001 ‘**’ 0.01 ‘*’ 0.05 ‘.’ 0.1 ‘ ’ 1

Residual standard error: 1.769 on 152 degrees of freedom
Multiple R-squared:  0.9805,    Adjusted R-squared:  0.9792 
F-statistic: 763.2 on 10 and 152 DF,  p-value: < 2.2e-16

```

Taking the linear variables into the model, for the training part we only get a neglectable improvement of the coefficient of determination.

Again we can use the functions from *library(leaps)* to search for “optimal” models.

---

	<i>Input</i>
--	--------------

```

library(leaps)
l1 <- leaps(x = fat[, c(6:15, 20:29)], y = fat$vol)

```

---

If we want to model the proportion of body fat directly, not the volume, we can construct the corresponding auxiliary variables.

---

	<i>Input</i>
--	--------------

```

fat$neckvolf <- fat$neckvol / fat$weightkg
fat$chestvolf <- fat$chestvol / fat$weightkg
fat$abdomenvolf <- fat$abdomenvol / fat$weightkg
fat$hipvolf <- fat$hipvol / fat$weightkg
fat$thighvolf <- fat$thighvol / fat$weightkg
fat$kneevolf <- fat$kneevol / fat$weightkg
fat$anklevolf <- fat$anklevol / fat$weightkg
fat$bicepvolf <- fat$bicepvol / fat$weightkg
fat$forearmvolf <- fat$forearmvol / fat$weightkg
fat$wristsvolf <- fat$wristsvol / fat$weightkg

```

---

We begin with a simple model. We use only one variable (`abdomenvolf`) from the group of variables describing the torso, and one variable (`wristvolf`) from the higher principal components. With this simple model, we achieve about the same quality as previously with the full set of variables.

```
Input
lm.volf <- lm(body.fat ~ abdomenvolf + wristvolf, data = fat, subset = fat$train)
summary(lm.volf)

Output
Call:
lm(formula = body.fat ~ abdomenvolf + wristvolf, data = fat,
subset = fat$train)

Residuals:
    Min      1Q  Median      3Q     Max 
-10.196 -3.021  0.100  3.161  9.164 

Coefficients:
            Estimate Std. Error t value Pr(>|t|)    
(Intercept) 0.2773356  6.7751858   0.041   0.967    
abdomenvolf 0.0022808  0.0001962  11.627 < 2e-16 ***  
wristvolf   -0.0331966  0.0054526  -6.088 8.16e-09 ***  
---
Signif. codes:  0 ‘***’ 0.001 ‘**’ 0.01 ‘*’ 0.05 ‘.’ 0.1 ‘ ’ 1

Residual standard error: 4.318 on 160 degrees of freedom
Multiple R-squared:  0.6762,    Adjusted R-squared:  0.6722 
F-statistic: 167.1 on 2 and 160 DF,  p-value: < 2.2e-16
```

So the recommendation would be: if you do not have a bath tube with you to determine the density, take a measurement on weight, height, abdomen and wrist.

a:ch04:16

Exercise 4.11	
*	Extend the variables by other volume related variables in the model given above. Do you gain precision?
**	Try to include the variable <code>age</code> in the model. How do you exactly include <code>age</code> in the model?
**	The function <code>mvr()</code> in <code>library(pls)</code> is available to perform a regression based on on principal components. Use this function for regression. What is the difference between this estimation and usual least squares regression?

For model construction, we used only the training part of the data. The quality of the model derived now can be checked using the evaluation part. This can be done using function `predict.lm()` which applies a model estimated with `lm()` to a new data set with analogous structure, for example

```
a:ch04:17
```

---

```
fat.eval <- fat[fat$train == FALSE, ]  
pred <- predict.lm(lm.volf, fat.eval, se.fit = TRUE)
```

<b>Exercise 4.12</b>	
*	Estimate the precision of the model using the evaluation part of the data.
*	Carry out a regression diagnostics of the model derived, using the evaluation part of the data.

## 4.8 High Dimensions

In small dimensions, we can analyse many problems completely. Higher dimension often require to design special strategies for analysis. Formal applications of standard methods soon meet their limitations.

Higher dimensions, such as dimension 10 through 100, are common in many application fields. But even problems in high dimensions occur day-to-day. We have to see that dimension is a question of modelling, not a mere question of the problem. Digital video (DV PAL) for example records images in a format of  $720 \times 576$ . A single image with three colours hence gives a vector in  $720 \times 576 \times 3 = 1244160$  dimensional space. A second of video has 25 of these frames. If we have to work with image data, it is our choice to see image processing as a problem in dimension  $d = 1244160$ , or as sequence of 1244160 (not independent!) observations in dimension  $d = 1$ .

Going from dimension  $d = 1244160$  to  $d = 1$  we shift information that is implicit in the dimensions to structural information. Of course, we have to pay for it: now the structural information has to be modelled.

As a remark: practical application takes a middle course. The image is decomposed into blocks, for example of size  $64 \times 64$ . Pixels within a block are handled simultaneously; the blocks are handled sequentially. You can see it the next time your digital TV has a glitch.

In high dimensional problems, the statistics is often not apparent at all. It is hidden in the hardware as “imbedded system”.

The figure 4.10 is an example from an [Rnews:Sawitzki:2002](#) analysis with R for a high dimensional data set of cDNA micro array data ([28]). A single observation in this data set consists of measurements on 4227 probes, each with 4 partial recordings for two colours (red and green) and two attributes (foreground fg and background bg), giving (*fg.green*, *fg.red*, *bg.green*, *bg.red*). This means one data point is of dimension  $4 \times 4227$ . This information is physically related to positions on the slide, the spots. The essential function, which is used for visualisation here is *image()*. It is used here to represent a variable *z* using a colour table against two co-ordinates *x*, *y* which give the position of a spot on a slide. The image shows one observation. The four channels for the partial recordings are put on different displays side by side, arranged by spot position. A unified yellow/blue colour table is chosen, which is recognizable even for most colour perception deficiencies.

In this kind part of the analysis, quality control was the main issue. A first information is to identify “dirty” spots, that is spots where the background intensity exceeds the foreground intensity. These are identified as red spots, and their distribution is represented in special

plots. The other critical information is on the adjustment of the channels. This information is represented in the bottom line. In this series of experiments, there should be a small fraction of differentially activated genes. So the distribution of the red intensity should be approximatively that of the green intensity. For this specific chip, we see that the estimated densities for the red and green dyes clearly differ, as shown in the bottom left plot. So this particular chip has a quality problem.

What we have illustrated here is an example of an analysis in dimension  $4 \times 4227$ . It is feasible, but of course this is not covered by standard methods.

The colours code the result of a pre-analysis - the red point mark problem areas on the cDNA-Chip. In this case, the pattern of the highlighted selection could identify a specific problem in the production process.

Theme oriented surveys about R packages, in particular for multivariate problems, are to be found in <<http://cran.at.r-project.org/src/contrib/Views/>>.

## 4.9 Statistical Summary

The analysis of multivariate data can only be touched in this context. Multivariate problems already occur implicitly for regression problems (see chapter 2). For simple regression problems, the multivariate aspects only referred to deterministic parameters. In the general case we have to analyse a multivariate statistical distribution. The introduction has to stop here, and more is reserved to subsequent lectures.

## 4.10 Literature and Additional References

```
Generated by Sweave from:
+ Source: /u/math/j40/cvsroot/lectures/src/SIntro/Rintro/Rnw/S04multiv.Rnw,v +
+ Revision: 1.34 +
+ Date: 2008/06/18 20:39:18 +
+ name: +
+ Author: j40 +
+Source: /u/math/j40/cvsroot/lectures/src/SIntro/Rintro/Rnw/S04multiv.Rnw,v +
+Revision : 1.34 +
+Date : 2008/06/18 20:39:18 +
+name : +
+Author : j40 +
```

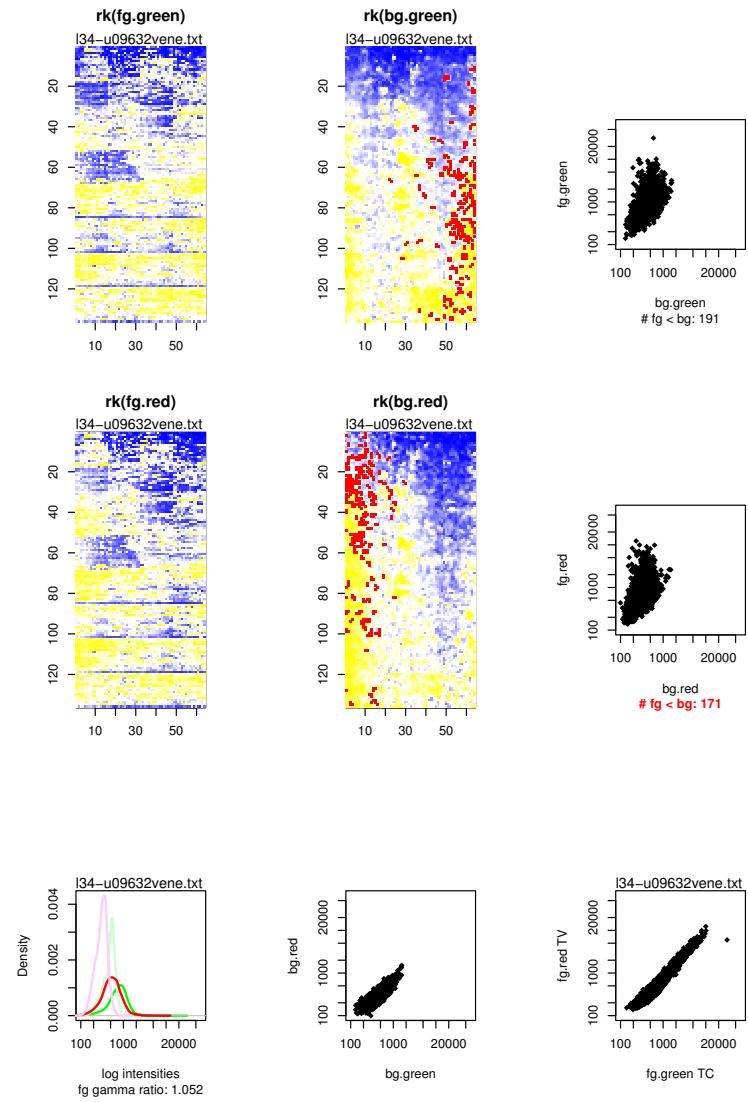
Figure 4.10 A  $4227 \times 4$  observation from a micro array experiment

fig:cDna



---

# R as a Programming Language and Environment

---

ch:0A

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

## A.1 Help and Information

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>apropos()</code>	localises by keyword.
<code>demos()</code>	executes demos for a topic area. <i>Example:</i> <code>demos(graphics)</code> <code>demos()</code> lists all topic areas which provide a demo

## 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 Search Paths	
<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>lsf.str("package:base")</code>
<code>find()</code>	localises by keyword. Does also find overlaid entries. <i>Syntax:</i> <code>find(what, mode = "any", numeric = FALSE, simple.words = TRUE)</code>
<code>apropos()</code>	localises by keyword. Does also find 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 Search Paths (cont.)	
<code>environment()</code>	current environments <i>Example:</i> <code>environment()</code>
<code>sys.parent()</code>	preceding environments <i>Example:</i> <code>sys.parent(1)</code>

### A.3 Customising

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.

**To Do:** An  
Introduction to R  
Sec. 10.8

<i>Some System Components with Global State</i>	
basic system	see <code>help(options)</code> .
random numbers	see Appendix A.23 (page A-36). <small>s:A.13 s:A.13</small>
basic graphics	see <code>help(par)</code> .
lattice graphics	see <code>help(lattice.options)</code> .

#### A.4 Basic Data Types

s:A.1

<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> 1.0 2 3.14E0
<i>complex</i>	complex, in Cartesian co-ordinates <i>Example:</i> 1.0+0i
<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> "T", 'klm'
<i>list</i>	general list structure List elements can be of different types. <i>Example:</i> list(1:10, "Hello")
<i>function</i>	R function <i>Example:</i> sin
<i>NULL</i>	special case: empty object <i>Example:</i> NULL

*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
<i>NA</i>	“not available”. Type: logical. NA 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

(cont.)→

<i>Special Constants</i> (cont.)	
<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>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

The object attributes and 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.

R <i>Inspection</i>	
<code>print()</code>	standard output
<code>structure()</code>	output, optional with attributes
<code>summary()</code>	standard output as summary, in particular for model fits
<code>plot()</code>	standard graphic output

### 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(&lt;object&gt;)</code>
<code>structure()</code>	shows the internal structure of an object. Attributes for the display can be passed as parameter. <i>Example:</i> <code>structure(1:6, dim = 2:3)</code> <i>Syntax:</i> <code>structure(&lt;object&gt;, ...)</code>
<code>class()</code>	object class. For newer objects, the class is stored as an attribute. In older S or R versions, 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(&lt;obj&gt;)</code> gives the <code>names</code> attribute of <code>&lt;obj&gt;</code> . <code>names(&lt;obj&gt;)&lt;-&lt;charvec&gt;</code> sets the <code>names</code> attribute. <i>Example:</i> <code>x&lt;-values</code> <code>names(x)&lt;- &lt;charvec&gt;</code>

**ToDo:**  
improve  
example  
for names

### 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(&lt;fun&gt;)</code> Shows specialised functions for <code>&lt;fun&gt;</code> , <code>methods(class = &lt;c&gt;)</code> the class specific functions for class <code>&lt;c&gt;</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

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 [2.6.5 page 2-43](#)).

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

`help(DateTimeClasses)`.

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

Compound Data Types	
<b>Vectors</b>	basic R data types
<b>Matrices</b>	vectors with two dimensional layout
<b>Arrays</b>	<p>vectors with higher dimensional layout</p> <p><code>dim()</code> defines a dimension attribute.</p> <p><i>Example:</i> <code>x &lt;- runif(100)</code>  <code>dim(x) &lt;- c(5, 5, 4)</code></p> <p><code>array()</code> generates a new vector with specified dimension structure.</p> <p><i>Example:</i> <code>z &lt;- array(0, c(4, 3, 2))</code></p> <p><code>rbind()</code> combines by rows.</p> <p><code>cbind()</code> combines by columns.</p>
<b>Factors</b>	<p>special case for categorical data</p> <p><code>factor()</code> converts a vector into a factor.</p> <p><i>See also</i> section <a href="#">2.2.1</a></p> <p><code>ordered()</code> converts a vector into a factor with ordered levels. This is a shortcut for <code>factor(x, ..., ordered = TRUE)</code></p> <p><code>levels()</code> returns the levels of a factor.</p> <p><i>Example:</i> <code>x &lt;- c("a", "b", "a", "c", "a")</code>  <code>xf &lt;- factor(x)</code>  <code>levels(xf)</code>  results in  <code>[ 1 ] "a" "b" "c"</code></p> <p><code>tapply()</code> applies a function separately for all levels of factors in a list.</p>

(cont.)→

<i>Compound Data Types</i> (cont.)	
<b>Lists</b>	<p>analogous to vectors, with elements of possibly different types</p> <p><code>list()</code> generates a list</p> <p>Syntax: <code>list(&lt;components&gt;)</code></p> <p><code>[[ ]]</code> access to components of a list by index</p> <p><code>&lt;list\$component&gt;</code> access by names</p> <p>Example: <code>l &lt;- list(name = "xyz", age = 22, fak = "math")</code>  <code>&gt;l[[2]]</code>  <code>22</code>  <code>&gt;l\$age</code>  <code>22</code></p>
<b>Data Frames</b>	<p><b>data frames</b> analogous to arrays resp. lists, with column-wise uniform type and uniform column length</p> <p><code>data.frame()</code> analogous to <code>list()</code>, but restrictions have to be satisfied.</p> <p><code>attach()</code> attaches a data to the current search list. For access to components the component name will be sufficient.</p> <p><code>detach()</code></p>

### A.9 Accessing Components

A.2

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 form of an expression rule.

<i>Accessing Components</i>	
$x[\langle\text{indices}\rangle]$	indicated components of $x$ <i>Example:</i> $x[1:3]$
$x[-\langle\text{indices}\rangle]$	$x$ omitting indicated components <i>Example:</i> $x[-3]$ $x$ omitting 3. components
$x[\langle\text{condition}\rangle]$	components of $x$ , for which the $\langle\text{condition}\rangle$ holds. <i>Example:</i> $x[x < 0.5]$

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

<i>R Index Access</i>	
<code>dim()</code>	gets or sets dimensions of an object. <i>Example:</i> $x <- 1:12$ ; <code>dim(x) &lt;- 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> refRexpl:ch01-MC06
<code>array()</code>	generated a possibly higher dimensional matrix. <i>Example:</i> <code>array(x, dim = length(x), dimnames = NULL)</code>

For manipulation of arrays and matrices, see Appendix A.9 (page A-12).

R Iterators	
<code>apply()</code>	applies a function to the rows or columns of a matrix. Syntax: <code>apply(x, MARGIN, FUNCTION, ...)</code> MARGIN = 1: rows, MARGIN = 2: columns See also refRexpl:ch01-MC06
<code>lapply()</code>	applies a function to the elements of a list. Syntax: <code>lapply(X, FUN, ...)</code>
<code>sapply()</code>	applies a function to the elements of a list, of a vector or a matrix. If possible, dimension names are carried over. Syntax: <code>sapply(X, FUN, ..., simplify = TRUE, USE.NAMES = TRUE)</code>
<code>mapply()</code>	applies a function to multiple list or vector arguments. Syntax: <code>mapply(FUN, ..., MoreArgs = NULL, simplify = TRUE, USE.NAMES = TRUE)</code>
<code>Vectorize()</code>	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. Syntax: <code>Vectorize(FUN, vectorize.args = arg.names, simplify = TRUE, USE.NAMES = TRUE)</code>
<code>tapply()</code>	applies a function to components of an object depending on a list of controlling factors.
<code>by()</code>	object-oriented variant of <code>tapply</code> Syntax: <code>by(data, INDICES, FUN, ...)</code>
<code>aggregate()</code>	calculates statistics for subsets. Syntax: <code>aggregate(x, ...)</code>
<code>replicate()</code>	evaluates an expression repeatedly (for example with generating random numbers for simulation). Syntax: <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. Syntax: <code>outer(vec1, vec2, FUNCTION, ...)</code>

A:datamamip

### A.10 Data Manipulation

<b>R Array Access</b>	
<code>cbind()</code> <code>rbind()</code>	combines columns resp. rows.
<code>split()</code>	splits a vector by factors.
<code>table()</code>	generates a table of counts. <code>rbind()</code> combines by rows. <code>cbind()</code> combines by columns.

<b>Transformations</b>	
<code>seq()</code>	generate a sequence.
<code>abbreviate()</code>	

<b>Transformations</b>	
<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

<b>Character String Transformations</b>	
<code>casefold()</code>	translates characters, in particular from upper to lower case or vice versa.
<code>tolower()</code>	translates to lower case.
<code>toupper()</code>	translates to upper case.
<code>chartr()</code>	translates characters in a character vector.
<code>substring()</code>	

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

<i>Vector Manipulation</i>	
<code>stack()</code>	concatenates multiple vectors from a data frame or list into a single vector and generates a factor indicating the source of each item. Syntax: <code>stack(x, ...)</code>
<code>unstack()</code>	splits a vector by an indicator variable, i.e. reverses the operation of <code>stack()</code> . Syntax: <code>unstack(x, ...)</code>
<code>split()</code>	splits a vector into the groups defined by a factor. Syntax: <code>split(x, f, drop = FALSE, ...)</code>
<code>unsplit()</code>	combines components to a vector, i.e. reverses <code>split()</code> . Syntax: <code>unsplit(value, f, drop = FALSE)</code>
<code>cut()</code>	converts a numeric to factor. <code>cut()</code> divides the range of a vector into intervals and creates a factor indicating the interval for each value. Syntax: <code>cut(x, ...)</code>

**ToDo:**

refer to re-

Shape  
 data  
 frames  
 in Hmisc  
 Harrell  
 4.2.7

<i>Transformations</i>	
<code>t()</code>	transposes rows and columns. Syntax: <code>t(x)</code>
<code>aperm()</code>	generalised permutation Syntax: <code>aperm(x, perm)</code>

(cont.)→

<b><i>Transformations</i></b> (cont.)	
	where <code>perm</code> is a permutation of the indices of <code>x</code> .
<code>split()</code>	splits a vector by a factor.
<code>unsplit()</code>	combines components to a vector.

### A.11 Operators

s:A.4

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>	
\$	components-selection <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%"&lt;-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 operand 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 write as a sequence with separation semicolons. Expression groups can be combined by {...}.

### A.12 Functions

s:A.5

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

<i>R Function Declarations</i>	
<b>Declarations</b>	<code>function (&lt;formal argument list&gt; ) &lt;expression&gt;</code> <i>Example:</i> <code>fak &lt;- function(n) prod(1:n)</code>
<b>Formal argument</b>	<code>&lt;argument name&gt;</code> <code>&lt;argument name&gt; = &lt;default value&gt;</code>
<b>Formal argument list</b>	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 &lt;- function (...)mean(c(...))</code>
<b>Function result</b>	<code>return &lt;value&gt;</code> stops function evaluation and returns value <code>&lt;value&gt;</code> as last expression in a function declaration: returns value
<b>Assignments</b>	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>&lt;&lt;-</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>&lt;Variable&gt;&lt;&lt;-&lt;value&gt;</code>

<i>R Function Call</i>	
<b>Function call</b>	<code>&lt;name&gt;(&lt;Supplied (actual) argument list&gt; )</code> <i>Example:</i> <code>fak(3)</code>
<b>Supplied argument list</b>	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 for a corresponding actual argument is missing for a formal argument. <i>Syntax:</i> <code>&lt;list of values&gt;</code> <code>&lt;argument name&gt; = &lt;values&gt;</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 [R:Gentleman+Thaka:2000](#) allows direct access to variables. Techniques to use this are described in [10].

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

```
"inc<-" <-function (x, value) x+value
x <- 10
inc(x)<- 3
x
```

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

### A.13 Debugging and Profiling

A:debug

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]>`.

`<return>` goes to the next instruction, if the function is under control of `debug`. Continuous with the expression evaluation if `browser` has been called directly.

`n` goes to the next instruction (also if `browser` has been called directly).

`cont` Continuous with the expression evaluation.

`c` Short for `cont`. Continuous the expression evaluation.

`where` Shows call nesting.

`Q` Stop execution and jumps back to base state.

<i>Debug Help</i>	
<code>browser()</code>	suspends execution and enters the browser mode. Syntax: <code>browser()</code>
<code>recover()</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>Q</code> you leave <code>recover()</code> Syntax: <code>recover()</code> Hint: 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. Syntax: <code>debug(&lt;function&gt;)</code>
<code>undebug()</code>	cancels debugger control for a function. Syntax: <code>undebug(&lt;function&gt;)</code>
<code>trace()</code>	marks a function for trace control. On subsequent calls to the function, the call is signaled together with its arguments. Syntax: <code>trace(&lt;function&gt;)</code>

(cont.)→

<b>Debug Help</b> (cont.)	
<b>untrace()</b>	cancels trace control for a function. Syntax: <code>untrace(&lt;function&gt;)</code>
<b>traceback()</b>	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>
<b>try()</b>	allows for user defined error handling. Syntax: <code>traceback(&lt;expression&gt;)</code>

To measure execution time in selected code ranges, R provides a “profiling”. However this is only available if R has been compiled with the appropriate options. The options innovated at compiling can be queried using `capabilities()`.

<b>Profiling Support</b>	
<code>system.time()</code>	returns the execution time of an expression. This function is available in all implementations. Syntax: <code>system.time(&lt;expr&gt;, &lt;gcFirst&gt;)</code>
<code>Rprof()</code>	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>
<code>Rprofmem()</code>	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>
<code>summaryRprof()</code>	summarizes the output of <code>Rprof()</code> and reports the timinb by function. Syntax: <code>summaryRprof(filename = "Rprof.out", chunksize = 5000, memory = c("none", "both", "tseries", "stats"), index = 2, diff = TRUE, exclude = NULL)</code>

**A.14 Control Structures**

A:ctl

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

(cont.)→

<i>R Control Structures</i> (cont.)	
<code>break</code>	terminates the current loop.
<code>next</code>	terminates the current loop cycle and advances to next.

**A.15 Administration and Customisation****s:A.7**

<i>objects()</i> <i>ls()</i>	lists the objects in the current search path.
<i>rm()</i>	removes indicated objects. Syntax: <i>rm(&lt;object list&gt;)</i>

### A.16 Input and Output to Data Streams

s:A.8

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("&lt;file name&gt; ")</code> <i>Example:</i> <code>source("cmnds.R")</code>
<code>sink()</code>	redirects output in the file specified. <i>Syntax:</i> <code>sink("&lt;file name&gt;")</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 = "&lt;dumpdata.R&gt;", append = FALSE)</code>

### A.17 External Data

**A.9Data**

To edit or enter data using a spreadsheet model, R provides `edit()` (previously called `data.entry()`).

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 `stataread` for Stata, `foreign` for SAS, Minitab and SPSS, `RODBC` for SQL. For more information, see the manual “Data Import/Export” ([20]).

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.18 (page A-27)).

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 us `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. <i>Syntax:</i> <code>save(&lt;names of the objects to be stored&gt;, file = &lt;file name&gt;, ...)</code>
<code>load()</code>	loads data from an external file. <i>Syntax:</i> <code>load(file = &lt;file name&gt;, ...)</code>
<code>data()</code>	loads data. <code>data()</code> can handle various file formats, if the access paths and file-names the follow R conventions. <i>Syntax:</i> <code>data(..., list = character(0), package = c(.packages(), .Autoloaded), lib.loc = .lib.loc)</code> <i>Example:</i> <code>data(crimes) # lädt die 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 which 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>	<p>reads data tables.</p> <p>Syntax: <code>read.table(file, header = FALSE, sep = "\t", ...)</code></p> <p>Examples: <code>read.table(&lt;file name&gt;, header = TRUE, sep = '\t')</code>  headers in row 1, row labels in column 1  <code>read.table(&lt;file name&gt;, header = TRUE, sep = '\t')</code>  now row number, headers in row 1,</p>
<code>write.table()</code>	<p>writes data table</p> <p>Syntax: <code>write.table(file, header = FALSE, sep = '\t', ...)</code></p> <p>Examples: <code>write.table(&lt;data frame&gt;, &lt;file name&gt;, header = TRUE, sep = '\t')</code>  headers in row 1, row labels in column 1  <code>write.table(&lt;data frame&gt;, &lt;file name&gt;, header = TRUE, sep = '\t')</code>  now row number, headers in row 1,</p>

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 [11]:

```
# date col in all numeric format yyyyymmdd
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.18 Libraries, Packages

s:A.9

External information can be stored in (text) files and Packages. In general, additional functions are provided as packages. Packages are loaded with  
`library()`.

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)
```

<i>Package Utilities</i>	
<code>library()</code>	loads an add-on package. Syntax: <code>library(package, ...)</code> See also section <a href="#">1.5.6</a> <code>subs:ch01-lib</code>
<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(&lt;name&gt;)</code>
<code>install.packages()</code>	installs add-on package in <code>&lt;lib&gt;</code> , downloading it from the archive <i>CRAN</i> or from specified archives Syntax: <code>install.packages(pkgs, lib, CRAN =getOption("CRAN"), ...)</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 = "&lt;anRpackage&gt;", list, ...)</code>

Detailed information for building R packages is in "Writing R Extensions" ([23]). [R:Ext](#)

### A.19 Linear Algebra

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>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

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

## A.20 Model Descriptions

s:A.10

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

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

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

R allows to specify models by giving the rules 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 levels of $M_1$ and $M_2$	
$\% \text{ in } \%$	$M_1 \% \text{ in } \% M_2$	modified tensor product	$a + b \% \text{ in } \% a$ corresponds to $a + a : b$
*	$M_1 * M_2$	“crossed”	$M_1 + M_2$ corresponds 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 “interactions” up to level $n$	
$I()$	$I(M)$	interpret $M$ . Terms in $M$ retain their original meaning; the result determines the model.	$Y \sim (1 + I(X^2))$ corresponds to $E(Y) = a + bX^2$

Table A.57 Wilkinson-Rogers Notation for Linear Models

ta:wrn

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

examples

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

$y \sim 0 + x$	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-fakcor 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>	update and re-fit, or change a model
<code>model.matrix()</code>	generate 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 <a href="#">ch:02</a>
<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.21 Graphic Functions

s:A.11

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 camera/object model. For information about lattice see `help(lattice)`, for a survey about the functions in lattice see `library(help = lattice)`. Informationen about the basic graphics system follow 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.21.1 High Level Graphics

<i>High Level Graphics</i>	
<code>plot()</code>	generic graphic output
<code>pairs()</code>	pair-wise scatterplots
<code>coplot()</code>	scatterplots, conditioned on covariates
<code>qqplot()</code>	QQ-Plot
<code>qqnorm()</code>	Gaussian QQ-Plot
<code>qqline()</code>	adds a line to a Gaussian QQ-Plot, passing through the first and third quantile.
<code>hist()</code>	histogram See also section 1.3.2, page 1-30
<code>boxplot()</code>	box&whisker plot
<code>dotplot()</code>	
<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

### A.21.2 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 which suppose that there is already a defined plot environment.

<b>Low Level Plotting</b>	
<code>points()</code>	generic function. Marks points at specified positions. Syntax: <code>points(x, ...)</code>
<code>lines()</code>	generic function. Joins points at specified positions. Syntax: <code>lines(x, ...)</code>
<code>abline</code>	adds a line (in several representations) to a plot. Syntax: <code>abline(a, b, ...)</code>
<code>polygon()</code>	adds polygon with specified vertices.
<code>axis()</code>	Adds axis.

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

<b>Interactions</b>	
<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>

### A.21.3 Annotations and Legends

The high-level-function 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)`.

This is complemented by low-level functions.

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

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()` 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.21.4 Graphic Parameters and Layout

<i>Parametrisations</i>	
<code>par()</code>	sets parameters for the basic graphics system. Syntax: see <code>help(par)</code> Example: <code>par(mfrow = c(m, n))</code> splits the graphic area in $m$ rows and $n$ columns, to be filled row-wise. <code>par(mfcol = c(m, n))</code> fills the area column by column.
<code>split.screen()</code>	splits the graphic area in parts. Syntax: <code>split.screen(figs, screen, erase = TRUE)</code> . If <code>figs</code> is a pair of two arguments, these will fix the number of rows and columns.. If <code>figs</code> is a matrix, each row gives the coordinates of a graphic area in relative coordinates [0 ... 1]. <code>split.screen()</code> can be nested.

(cont.)→

<i>Parametrisations</i> (cont.)	
<code>screen()</code>	selects graphic area for the next graphical output. Syntax: <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.22 Basic Statistical Functions

s:A.12

<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> <i>min()</i>	maximum, minimum. See also <i>pmax</i> , <i>pmin</i> .
<i>range()</i>	minimum and maximum
<i>cummax()</i> <i>cummin()</i>	cumulated maximum, minimum
<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> <i>rev()</i>	sorting
<i>order()</i>	returns a permutation for sorting.
<i>rev()</i>	reverse sorting
<i>range()</i>	ranks in sample

### A.23 Distributions, Random Numbers, Densities...

s:A.13

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(&lt;name&gt;)</code> sets a basic random number generator.</p> <p>Syntax: <code>RNGkind()</code>  <code>RNGkind(&lt;name&gt;)</code></p> <p>Example: <code>RNGkind("Wichmann-Hill")</code>  <code>RNGkind("Marsaglia-Multicarry")</code>  <code>RNGkind("Super-Duper")</code></p>

(cont.)→

<b>R Random Numbers</b> (cont.)	
<b>sample()</b>	<p>draws a sample from the values given in vector <i>x</i>, with or without replacement (controlled by the value of <i>replace</i>).</p> <p>Size is by default the length of <i>x</i>.</p> <p>Optionally, <i>prob</i> can be a vector of probabilities for the values of <i>x</i>.</p> <p><i>Syntax:</i>    <code>sample(x, size, replace = FALSE, prob)</code></p> <p><i>Example:</i>    Random permutation:  <code>sample(x)</code></p> <pre>val&lt;-c("H", "T") prob&lt;-c(0.3, 0.7) sample(val, 10, replace = TRUE, prob)</pre>

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

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

With

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

the state of the random number generator is restored when necessary.

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

<i>xxxxx</i>	generates random numbers
<i>dxxxxx</i>	density or probability
<i>pxxxx</i>	distribution function
<i>qxxxx</i>	quantiles

Example:

<code>x&lt;-runif(100)</code>	generates 100 random variable with $U(0, 1)$ distribution.
<code>qf(0.95, 10, 2)</code>	calculates the 95% quantile of the $F(10, 2)$ distribution.

<i>Distributions</i>	<i>Short Name</i>	<i>Parameter and Default-values</i>
Beta	<code>beta</code>	<code>shape1, shape2, ncp = 0</code>

(cont.)→

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

### A.24 Computing on the Language

The language elements of R are objects, as are data or functions. They can be read and changed like these. Chapter 6 of the “R Language Definition” [21] gives details for computing on the language. See also section 2.1.5

#### *Function objects*

of [21].

<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 &lt;- expression(integral(fun, lims))</code> <i>See also</i> <a href="#"> 1.3.1: mathematical typesetting in plot annotations</a>
<code>substitute()</code>	R expression with evaluation of all defined terms.
<code>bquote()</code>	R expression with selective evaluation. Terms in .() are evaluated. <i>Examples:</i> <code>n&lt;-10; bquote( n^2 == .(n*n))</code>

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

```
$Source: /u/math/j40/cvsroot/lectures/src/SIntro/Rintro/Rnw/S0A1.Rnw.tex,v $
$Revision: 1.36 $
$Date: 2008/06/18 20:39:19 $
$Name:  $
$Author: j40 $
```



## References

- [1] Y. Amit, U. Grenander, and M. Piccioni. Structural image restoration through deformable templates. *Journal of the American Statistical Association*, 86:376–387, 1991.
- [2] Richard A. Becker, John M. Chambers, and Allan R. Wilks. *The New S Language*. Chapman & Hall, London, 1988.
- [3] D. Brandis, V. Storch, and M. Türkay. Morphology and function of the reproductive system in the freshwater crab genus potamon. *Journal of Morphology*, 239:157–166, 1999.
- [4] D. Brandis, V. Storch, and M. Türkay. Taxonomy and zoogeography of the freshwater crabs of europe, north africa, and the middle east (crustacea: Decapoda: Potamidae). *Senckenbergiana Biologica*, 80:5–56, 2000.
- [5] John M. Chambers. *Programming with Data*. Springer, New York, 1998. ISBN 0-387-98503-4.
- [6] John M. Chambers and Trevor J. Hastie. *Statistical Models in S*. Chapman & Hall, London, 1992.
- [7] William S. Cleveland. *Visualizing Data*. AT&T Bell Laboratories, Murray Hill, 1993.
- [8] George W. Furnas and Andreas Buja. Prosection views: dimensional inference through sections and projections. *J. Comput. Graph. Statist.*, 3(4):323–385, 1994.
- [9] Peter Gänßler and Winfried Stute. *Wahrscheinlichkeitstheorie*. Springer, 1977.
- [10] Robert Gentleman and Ross Ihaka. Lexical scope and statistical computing. *Journal of Computational and Graphical Statistics*, 9:491–508, 2000.
- [11] Gabor Grothendieck and Thomas Petzoldt. R help desk: Date and time classes in R. *R News*, 4(1):29–32, June 2004.
- [12] Rob J. Hyndman, David M. Bashtannyk, and Gary K. Grunwald. Estimating and visualizing conditional densities. *J. Comput. Graph. Statist.*, 5(4):315–336, 1996.
- [13] Alfred Inselberg, Tuval Chomut, and Mordechai Reif. Convexity algorithms in parallel coordinates. *J. Assoc. Comput. Mach.*, 34(4):765–801, 1987.
- [14] N.L. Johnson and S. Kotz. *Discrete Distributions*. Wiley, New York, 1970.
- [15] Bent Jørgensen. *The Theory of Linear Models*. Chapman & Hall, New York-London, 1993.
- [16] Paul Massart. The tight constant in the Dvoretzky-Kiefer-Wolfowitz inequality. *The Annals of Probability*, 18(3):1269–1283, July 1990.
- [17] R. G. Miller. *Simultaneous Statistical Inference*. Springer, New York, 1981.
- [18] Paul Murrell. *R Graphics*. Chapman & Hall/CRC, Boca Raton, Fla. [u.a.], 2006.
- [19] R Development Core Team. An introduction to R. Technical report, R Project, 2007.
- [20] R Development Core Team. R data import/export. Technical report, R Project, 2007.
- [21] R Development Core Team. The R language definition. Technical report, R Project, 2007.
- [22] R Development Core Team. The R reference index. Technical report, R Project, 2007.
- [23] R Development Core Team. Writing R extensions. Technical report, R Project, 2008.
- [24] C. Radhakrishna Rao. *Linear Statistical Inference and Its Applications*. Wiley, 2 edition, 1973.
- [25] G. M. Reaven and R. G. Miller. An attempt to define the nature of chemical diabetes using a multidimensional analysis. *Diabetologia*, 16:17–24, 1979.
- [26] Günther Sawitzki. Numerical reliability of data analysis systems. *Computational Statistics & Data Analysis*, 18(2):269–286, 1994.

## Bib-2

## Bibliography

**gs94nrel2**

- [27] Günther Sawitzki. Report on the numerical reliability of data analysis systems. *Computational Statistics & Data Analysis*, 18(2):289 – 301, 1994.

**news:Sawitzki:2002**

- [28] Günther Sawitzki. Quality control and early diagnostics for cDNA microarrays. *R News*, 2(1):6–10, March 2002.

**006Graphics-of-Lar**

- [29] Antony Unwin, Martin Theus, and Heike Hofmann. *Graphics of Large Datasets: Visualizing a Million*. Springer, 2006.

**vr00pis**

- [30] William N. Venables and Brian D. Ripley. *S Programming*. Springer, 2000. ISBN 0-387-98966-8.

**vr02mass**

- [31] William N. Venables and Brian D. Ripley. *Modern Applied Statistics with S*. Springer, Heidelberg, 4 edition, 2002.

---

# Functions by Topic

---

**Topic NA**

`is.type`, 2-43, A-4  
`is.na`, 2-43, A-5

**Topic algebra**

`eigen`, A-28  
`matrix`, A-11  
`qr`, A-28  
`solve`, A-28  
`svd`, A-28

**Topic aplot**

`abline`, 1-14  
`axis`, A-32  
`contour`, 4-2, 4-3, A-31  
`coplot`, 4-18, A-31  
`image`, 4-2, 4-3, 4-5, 4-55, A-31  
`legend`, 1-51, A-33  
`lines`, A-32  
`mtext`, 1-47, A-33  
`points`, A-32  
`polygon`, A-32  
`rug`, 1-17  
`screen`, A-34  
`split.screen`, A-33  
`title`, A-33

**Topic arith**

`cummax`, A-35  
`cummin`, A-35  
`cumprod`, A-35  
`cumsum`, A-35  
`max`, A-35  
`min`, A-35  
`prod`, A-35  
`range`, A-35  
`sort`, 1-14, A-35  
`sum`, A-35

**Topic array**

`aggregate`, A-12  
`aperm`, A-14  
`apply`, 1-25, A-12  
`array`, A-9, A-11  
`cbind`, A-9, A-13

`determinant`, A-28

`dim`, A-9, A-11  
`dimnames`, A-11  
`eigen`, A-28  
`expand.grid`, A-14  
`matrix`, A-11  
`merge`, A-14  
`ncol`, A-11  
`nrow`, A-11  
`outer`, 1-25, A-12  
`qr`, A-28  
`rbind`, A-9, A-13  
`svd`, A-28  
`t`, A-14  
`var`, 1-36, A-35

**Topic attribute**

`attr`, 2-43  
`attributes`, A-7  
`length`, 1-14, A-7, A-35  
`mode`, 2-43, A-4, A-7  
`names`, A-7  
`storage.mode`, 2-43, A-7  
`structure`, A-6, A-7  
`typeof`, 2-43, A-4, A-7

**Topic category**

`by`, A-12  
`cut`, A-14  
`factor`, 2-5, A-9  
`levels`, A-9  
`ordered`, A-9  
`split`, A-13–A-15  
`table`, 1-21, A-13, A-14  
`tapply`, A-9, A-12  
`unsplit`, A-14, A-15

**Topic character**

`abbreviate`, A-13  
`casefold`, A-13  
`chartr`, A-13  
`pmatch`, A-13  
`substring`, A-13  
`tolower`, A-13

## Index-2

`toupper`, A-13  
Topic **chron**  
  `rep`, 1-9, Private-7  
Topic **classes**  
  `class`, 2-43, A-7  
  `data.frame`, A-10  
  `inherits`, 2-43  
  `unclass`, 2-43  
Topic **color**  
  `col2rgb`, 1-53  
  `colours`, 1-53  
  `rgb`, 1-53  
Topic **connection**  
  `deparse`, A-39  
  `dump`, A-24  
  `parse`, 1-54  
  `read.fwf`, A-26  
  `read.table`, A-26  
  `scan`, A-26  
  `sink`, A-24  
  `source`, 1-56, 1-57, A-24  
  `write`, A-24  
Topic **datasets**  
  `data`, 1-57, 2-41, A-8, A-25, A-27  
Topic **data**  
  `apropos`, A-1, A-2  
  `attach`, A-10  
  `bquote`, 1-26, A-33, A-39  
  `deparse`, A-39  
  `detach`, A-10, A-27  
  `environment`, 1-55, A-2, A-18  
  `eval`, 1-55, A-39  
  `find`, A-2  
  `library`, A-8, A-27  
  `require`, A-27  
  `search`, 1-55, A-1, A-2, A-8  
  `searchpaths`, A-2  
  `substitute`, A-39  
  `sys.parent`, A-2  
Topic **debugging**  
  , A-19  
  `browser`, A-19  
  `debug`, A-19  
  `recover`, A-19  
  `traceback`, A-19, A-20  
  `untrace`, A-20  
Topic **design**  
  `contrasts`, A-30  
Topic **device**  
  `split.screen`, A-33

## Functions by Topic

Topic **diplot**  
  `identify`, 4-46  
Topic **distribution**  
  `.Random.seed`, 1-8, A-36  
  `chisq.test`, 1-31  
  `density`, 1-17  
  `hist`, 1-17, 1-20, 4-5, A-31  
  `qqnorm`, 3-6, 3-7, 4-5, A-31  
  `qqplot`, 3-7, A-31  
  `RNGkind`, A-36  
  `sample`, A-37  
  `set.seed`, A-36  
  `Uniform`, 1-4, Private-1  
Topic **documentation**  
  `apropos`, A-1, A-2  
  `args`, A-1  
  `citation`, 1-57  
  `data`, 1-57, 2-41, A-8, A-25, A-27  
  `demo`, A-1  
  `example`, A-1  
  `find`, A-2  
  `help`, A-1, A-8  
  `help.search`, A-1  
  `help.start`, 1-3, A-1  
  `package.manager`, A-27  
  `str`, A-7  
Topic **dplot**  
  `cloud`, 4-3  
  `col2rgb`, 1-53  
  `colours`, 1-53  
  `densityplot`, 4-5  
  `expression`, 1-15, A-33, A-39  
  `hist`, 1-17, 1-20, 4-5, A-31  
  `lattice.options`, A-8  
  `par`, A-8  
  `parallel`, 4-5  
  `persp`, 4-2, 4-3, 4-5, A-31  
  `qq`, 4-5  
  `split.screen`, A-33  
  `trellis.par.set`, A-8  
  `wireframe`, 4-3  
Topic **environment**  
  `apropos`, A-1, A-2  
  `browser`, A-19  
  `debug`, A-19  
  `find`, A-2  
  `ls`, 1-55, A-2, A-8, A-23  
  `objects`, A-2, A-23  
  `options`, A-8  
  `par`, A-8

**q**, 1-3  
**rm**, A-23  
**undebug**, A-19  
**Topic error**  
  **options**, A-8  
**Topic file**  
  **data.entry**, A-25  
  **dir**, A-8  
  **dump**, A-24  
  **load**, 2-41, A-25  
  **package.skeleton**, 1-57, 1-58, A-27  
  **read.fwf**, A-26  
  **read.table**, A-26  
  **save**, 1-57, 2-41, A-25  
  **scan**, A-26  
  **sink**, A-24  
  **source**, 1-56, 1-57, A-24  
  **system**, A-8  
  **write**, A-24  
  **write.table**, A-26  
**Topic hplot**  
  **barchart**, 4-5  
  **barplot**, 4-5  
  **boxcox**, 2-36  
  **boxplot**, 1-38, 4-5, A-31  
  **bwplot**, 4-5  
  **cloud**, 4-3, 4-5, 4-13  
  **contourplot**, 4-5  
  **coplot**, 4-18  
  **curve**, 1-23, A-31  
  **dotchart**, 4-5  
  **dotplot**, 4-5, A-31  
  **hist**, 1-17, 1-20, 4-5, A-31  
  **histogram**, 4-5  
  **pairs**, 4-5, 4-8, 4-12, 4-44, A-31  
  **persp**, 4-2, 4-3, 4-5, A-31  
  **plot**, 1-14, 1-15, 2-32, 2-43, 4-5, 4-25, 4-46, A-6, A-31  
  **qqmath**, 4-5  
  **qqnorm**, 3-6, 3-7, 4-5, A-31  
  **qqplot**, 3-7, 4-5, A-31  
  **splom**, 4-5  
  **stripchart**, 4-5  
  **stripplot**, 4-5, 4-26  
  **wireframe**, 4-3, 4-5  
  **xyplot**, 4-5  
**Topic htest**  
  **chisq.test**, 1-31  
  **glht**, 2-31  
  **ks.test**, 1-29

**power.prop.test**, 3-27  
**power.t.test**, 3-24  
**prop.test**, 3-26  
**t.test**, 3-13  
**wilcox.test**, 3-15, 3-16  
**Topic interface**  
  **system**, A-8  
**Topic iplot**  
  **locator**, A-32  
  **par**, A-8  
**Topic iteration**  
  **apply**, 1-25, A-12  
  **by**, A-12  
  **lapply**, A-12  
  **mapply**, A-12  
  **replicate**, A-12  
  **sapply**, A-12  
  **tapply**, A-9, A-12  
  **Vectorize**, A-12  
**Topic list**  
  **lapply**, A-12  
  **list**, A-10  
  **replicate**, A-12  
  **sapply**, A-12  
  **Vectorize**, A-12  
**Topic loess**  
  **loess**, 2-37  
**Topic logic**  
  **duplicated**, A-13  
  **is.<type>**, 2-43, A-4  
  **is.na**, A-5  
  **is.nan**, A-5  
  **match**, A-13  
  **unique**, A-13  
**Topic manip**  
  **as.<type>**, 2-43, A-4  
  **c**, 1-9, Private-6  
  **cbind**, A-9, A-13  
  **cut**, A-14  
  **deparse**, A-39  
  **dimnames**, A-11  
  **duplicated**, A-13  
  **is.<type>**, 2-43, A-4  
  **list**, A-10  
  **mapply**, A-12  
  **match**, A-13  
  **merge**, A-14  
  **order**, A-35  
  **rbind**, A-9, A-13  
  **rep**, 1-9, Private-7

## Index-4

`reshape`, A-14  
`rev`, A-35  
`seq`, 1-9, A-13, Private-7  
`sort`, 1-14, A-35  
`split`, A-14  
`stack`, 2-41, A-14  
`structure`, A-6, A-7  
`unique`, A-13  
`unsplit`, A-14  
`unstack`, A-14  
`Vectorize`, A-12  
Topic **math**  
`integrate`, 4-34  
`is.inf`, 2-43, A-5  
`is.nan`, 2-43, A-5  
Topic **methods**  
`class`, 2-43, A-7  
`data.frame`, A-10  
`inherits`, 2-43  
`methods`, 2-44, A-8  
`predict`, 2-27, 2-45  
`summary`, 1-37, A-6  
`unclass`, 2-43  
`UseMethod`, 2-43, 2-44  
Topic **misc**  
`citations`, 1-57  
`wilcox_test`, 3-15  
Topic **models**  
`anova`, 2-22, 2-31, A-30  
`aov`, 2-23  
`boxcox`, 2-36  
`coef`, 2-44  
`confint`, 2-45  
`effects`, 2-44  
`expand.grid`, A-14  
`fitted`, 2-45  
`formula`, A-30  
`glm`, A-30  
`model.matrix`, 2-18, 2-45, A-30  
`nls`, A-30  
`residuals`, 2-44  
`stdres`, 2-44  
`studres`, 2-45  
`terms`, A-30  
`update`, A-30  
`vcov`, 2-45  
Topic **multivariate**  
`lda`, 4-27  
`mvr`, 4-54  
`prcomp`, 4-30, 4-51

## Functions by Topic

`var`, 1-36, A-35  
Topic **nonlinear**  
`nlm`, A-30  
`nls`, A-30  
`vcov`, 2-45  
Topic **optimize**  
`nlm`, A-30  
Topic **print**  
`ls.str`, A-2  
`options`, A-8  
`print`, 1-55, 2-43, 4-4, 4-13, A-6  
`str`, A-7  
`write.table`, A-26  
Topic **programming**  
`bquote`, 1-26, A-33, A-39  
`browser`, A-19  
`debug`, A-19  
`deparse`, 1-55, A-39  
`environment`, 1-55, A-2, A-18  
`eval`, 1-55, A-39  
`expression`, 1-15, A-33, A-39  
`missing`, 1-49, A-17  
`parse`, 1-54, A-39  
`recover`, A-19  
`source`, 1-56, 1-57, A-24  
`substitute`, 1-55, A-39  
`sys.parent`, A-2  
`trace`, A-19  
 `traceback`, A-20  
`try`, A-20  
`udebug`, A-19  
`untrace`, A-20  
Topic **regression**  
`anova`, 2-22, 2-31, A-30  
`aov`, 2-23  
`boxcox`, 2-36  
`coef`, 2-44  
`contrasts`, A-30  
`effects`, 2-44  
`fitted.`, 2-45  
`glm`, A-30  
`influence`, 2-45  
`leaps`, 4-50  
`lm`, 2-6, 2-12, 2-23, 2-44, 4-47, 4-54, A-30  
`mvr`, 4-54  
`nls`, A-30  
`predict.lm`, 2-27, 4-54  
`regsubsets`, 4-49  
`residuals`, 2-44  
Topic **robust**

**median**, A-35  
Topic **smooth**  
  **density**, 1-17, A-30  
  **loess**, 2-37  
Topic **sysdata**  
  **.Random.seed**, 1-8, A-36  
  **colours**, 1-53  
  **RNGkind**, A-36  
  **set.seed**, A-36  
Topic **univar**  
  **max**, A-35  
  **mean**, 1-36, A-35  
  **median**, A-35  
  **min**, A-35  
  **order**, A-35  
  **quantile**, 1-37, A-35  
  **range**, A-35  
  **sd**, 1-36  
  **sort**, 1-14, A-35  
  **var**, 1-36, A-35  
Topic **utilities**  
  **data.entry**, A-25  
  **demo**, A-1  
  **edit**, A-25  
  **example**, A-1  
  **getwd**, A-8  
  **install.packages**, 1-56, A-27  
  **integrate**, 4-34  
  **ls.str**, A-2  
  **mapply**, A-12  
  **package.skeleton**, 1-57, 1-58, A-27  
  **Rprof**, A-20  
  **Rprofmem**, A-20  
  **setwd**, A-8  
  **str**, A-7  
  **summaryRprof**, A-20  
  **Sweave**, 1-56  
  **system**, A-8  
  **system.time**, A-20  
  **update.packages**, 1-56  
  **Vectorize**, A-12



---

# Function and Variable Index

---

.Random.seed, 1-5, *Private-2*  
[, *Private-2*

aa1 itfun (*alias1*), *Private-1*  
abbreviate, A-13  
abline, 1-14  
add1, 2-23  
aggregate, A-12  
anova, 2-13, 2-22, 2-31, A-30  
anova.lm, 2-15  
anscombe, 2-15  
aov, 2-12, 2-13, 2-15, 2-23  
aperm, A-14  
apply, 1-25, A-12  
apropos, A-1, A-2  
args, A-1  
array, A-9, A-11  
as.<type>, 2-43, A-4  
as.data.frame, 2-12, 2-38  
attach, A-10  
attitude, 2-15  
attr, 2-43  
attributes, A-7  
axis, A-32

barchart, 4-5  
barplot, 4-5  
boxcox, 2-36  
boxplot, 1-38, 4-5, A-31  
bquote, 1-26, A-33, A-39  
browser, A-19  
bwplot, 4-5  
by, A-12

c, 1-9, *Private-6*  
casefold, A-13  
cbind, A-9, A-13  
chartr, A-13  
chisq.test, 1-31  
citations, 1-57

class, 2-13, 2-43, A-7  
cloud, 4-3, 4-5, 4-13  
co.intervals (*coplot*), 4-18  
coef, 2-15, 2-44  
coefficients, 2-23  
col2rgb, 1-53  
colours, 1-53  
confint, 2-15, 2-45  
contour, 4-2, 4-3, A-31  
contourplot, 4-5  
contrasts, A-30  
coplot, 4-18, 4-18, A-31  
cummax, A-35  
cummin, A-35  
cumprod, A-35  
cumsum, A-35  
curve, 1-23, A-31  
cut, A-14

data, 1-57, 2-41, A-8, A-25, A-27  
data.entry, A-25  
data.frame, A-10  
data.matrix, 4-9  
debug, A-19  
demo, A-1  
density, 1-17  
densityplot, 4-5  
deparse, 1-55, A-39  
detach, A-10, A-27  
determinant, A-28  
dim, A-9, A-11  
dimnames, A-11  
dir, A-8  
dotchart, 4-5  
dotplot, 4-5, A-31  
drop1, 2-23  
dump, A-24  
dunif (*Uniform*), 1-4, *Private-1*  
duplicated, A-13

## Index-8

**edit**, A-25  
**effects**, 2-14, 2-15, 2-23, 2-44  
**eigen**, A-28  
**environment**, 1-55, A-2, A-18  
**eval**, 1-55, A-39  
**example**, A-1  
**expand.grid**, A-14  
**expression**, 1-15, A-33, A-39  
  
**factor**, 2-5, 4-19, A-9  
**find**, A-2  
**fitted**, 2-15, 2-45  
**fitted.values**, 2-23  
**formula**, 2-12, 2-13, 2-38, A-30  
**freeny**, 2-15  
**fun1**, Private-1  
**function**, 4-19  
  
**getwd**, A-8  
**glht**, 2-31  
**glm**, 2-15, A-30  
  
**help**, A-1, A-8  
**help.search**, A-1  
**help.start**, 1-3, A-1  
**hist**, 1-17, 1-20, 4-5, A-31  
**histogram**, 4-5  
  
**identify**, 4-46  
**image**, 4-2, 4-3, 4-5, 4-55, A-31  
**influence**, 2-45  
**inherits**, 2-43  
**install.packages**, 1-56, A-27  
**integrate**, 4-34  
**is.<type>**, 2-43, A-4  
**is.inf**, 2-43, A-5  
**is.na**, 2-43, A-5  
**is.nan**, 2-43, A-5  
  
**kruskal.test**, 3-18  
**ks.test**, 1-29  
  
**lapply**, A-12  
**lattice.options**, A-8  
**lda**, 4-27  
**leaps**, 4-50  
**legend**, 1-51, A-33  
**length**, 1-14, A-7, A-35

## Function and Variable Index

**levels**, A-9  
**library**, A-8, A-27  
**LifeCycleSavings**, 2-15  
**lines**, A-32  
**list**, A-10  
**lm**, 2-6, 2-12, 2-23, 2-44, 4-47, 4-54, A-30  
**lm.fit**, 2-13, 2-15  
**lm.influence**, 2-15  
**lm.wfit**, 2-15  
**load**, 2-41, A-25  
**locator**, A-32  
**loess**, 2-37  
**loess.control**, 2-38, 2-39  
**longley**, 2-15  
**lowess**, 2-39  
**ls**, 1-55, A-2, A-8, A-23  
**ls.str**, A-2  
  
**mapply**, A-12  
**match**, A-13  
**matrix**, 4-20, A-11  
**max**, A-35  
**mean**, 1-36, A-35  
**median**, A-35  
**merge**, A-14  
**methods**, 2-44, A-8  
**min**, A-35  
**missing**, 1-49, A-17  
**mode**, 2-43, A-4, A-7  
**model.frame**, 3-13, 3-17  
**model.matrix**, 2-13, 2-18, 2-45, A-30  
**model.matrix.default**, 2-13  
**model.offset**, 2-13  
**mtext**, 1-47, 4-20, A-33  
**mvr**, 4-54  
  
**NA**, 3-8  
**NA (*is.nan*)**, 2-43  
**na.exclude**, 2-12  
**na.fail**, 2-12  
**na.omit**, 2-12  
**names**, A-7  
**ncol**, A-11  
**nlm**, A-30  
**nls**, A-30  
**nrow**, A-11  
  
**objects**, A-2, A-23  
**offset**, 2-13

## Function and Variable Index

## Index-9

**options**, 2-12, A-8  
**order**, A-35  
**ordered**, A-9  
**outer**, 1-25, A-12  
  
**package.manager**, A-27  
**package.skeleton**, 1-57, 1-58, A-27  
**pairs**, 4-5, 4-8, 4-8, 4-12, 4-20, 4-44, A-31  
**panel.smooth**, 4-20  
**par**, 4-20, A-8  
**parallel**, 4-5  
**parse**, 1-54, A-39  
**persp**, 4-2, 4-3, 4-5, A-31  
**plot**, 1-14, 1-15, 2-32, 2-43, 4-5, 4-25, 4-46,  
 A-6, A-31  
**pmatch**, A-13  
**points**, 4-20, A-32  
**polygon**, A-32  
**power.prop.test**, 3-27  
**power.t.test**, 3-24  
**ppoints**, 3-8  
**prcomp**, 4-30, 4-51  
**predict**, 2-15, 2-27, 2-45  
**predict.lm**, 2-14, 2-15, 2-27, 4-54  
**predict.loess**, 2-39  
**print**, 1-55, 2-43, 4-4, 4-13, A-6  
**print.anova** (*anova*), 2-22  
**print.lm** (*lm*), 2-12  
**prod**, A-35  
**prop.test**, 3-14, 3-26  
**psignrank**, 3-18  
**punif** (*Uniform*), 1-4, Private-1  
**pwilcox**, 3-18  
  
**q**, 1-3  
**qq**, 4-5  
**qqline** (*qqnorm*), 3-7, A-31  
**qqmath**, 4-5  
**qqnorm**, 3-6, 3-7, 4-5, A-31  
**qqplot**, 3-7, 4-5, A-31  
**qqplot** (*qqnorm*), 3-7  
**qr**, A-28  
**quantile**, 1-37, A-35  
**qunif** (*Uniform*), 1-4, Private-1  
  
**range**, 4-20, A-35  
**rbind**, A-9, A-13  
**read.fwf**, A-26  
**read.table**, A-26  
  
**recover**, A-19  
**regsubsets**, 4-49  
**rep**, 1-9, *Private*-7  
**replicate**, A-12  
**require**, A-27  
**reshape**, A-14  
**residuals**, 2-15, 2-23, 2-44  
**rev**, A-35  
**rgb**, 1-53  
**rm**, A-23  
**RNGkind**, A-36  
**rnorm**, 1-5, *Private*-2  
**Rprof**, A-20  
**Rprofmem**, A-20  
**rug**, 1-17  
**runif**, 1-4  
**runif** (*Uniform*), 1-4, Private-1  
  
**sample**, A-37  
**sapply**, A-12  
**save**, 1-57, 2-41, A-25  
**scan**, A-26  
**screen**, A-34  
**sd**, 1-36  
**search**, 1-55, A-1, A-2, A-8  
**searchpaths**, A-2  
**seq**, 1-9, A-13, *Private*-7  
**set.seed**, A-36  
**setwd**, A-8  
**sink**, A-24  
**solve**, A-28  
**some fun**, Private-1  
**sort**, 1-14, A-35  
**source**, 1-56, 1-57, A-24  
**split**, A-13-A-15  
**split.screen**, A-33  
**splom**, 4-5  
**stack**, 2-41, A-14  
**stackloss**, 2-15  
**stdres**, 2-44  
**storage.mode**, 2-43, A-7  
**str**, A-7  
**stripchart**, 4-5  
**stripplot**, 4-5, 4-26  
**structure**, A-6, A-7  
**studres**, 2-45  
**substitute**, 1-55, A-39  
**substring**, A-13  
**sum**, A-35  
**summary**, 1-37, 2-23, A-6

## Index-10

**summary.lm**, 2-15  
**summaryRprof**, A-20  
**svd**, A-28  
**Sweave**, 1-56  
**swiss**, 2-15  
**sys.parent**, A-2  
**system**, A-8  
**system.time**, A-20

**t**, A-14  
**t.test**, 3-13, 3-18  
**table**, 1-21, A-13, A-14  
**tapply**, A-9, A-12  
**terms**, 2-14, A-30  
**title**, 4-20, A-33  
**tolower**, A-13  
**toupper**, A-13  
**trace**, A-19  
**traceback**, A-20  
**trellis.par.set**, A-8  
**try**, A-20  
**ts.intersect**, 2-14  
**typeof**, 2-43, A-4, A-7

**unclass**, 2-43  
**undebug**, A-19  
**Uniform**, 1-4, **Private-1**  
**unique**, A-13  
**unsplit**, A-14, A-15  
**unstack**, A-14  
**untrace**, A-20  
**update**, A-30  
**update.packages**, 1-56  
**UseMethod**, 2-43, 2-44

**var**, 1-36, A-35  
**vcov**, 2-15, 2-45  
**Vectorize**, A-12

**wilcox.exact**, 3-18  
**wilcox.test**, 3-15, **3-16**  
**wilcox\_test**, 3-15, 3-18  
**wireframe**, 4-3, 4-5  
**write**, A-24  
**write.table**, A-26

**xyplot**, 4-5

## Function and Variable Index

---

# Index

---

*PP* plot, **1-41**

*QQ* plot, **1-41**

Note

- A.24: ? show plot .RB, 1-10
- A.24: ?consider simple exercises for random numbers – Antony, 1-8
- A.24: ?get damaged/destroyed beyond repair, 1-10

ToCheck

- A.24: micronuclei still available?, 2-24

ToDo

- A.24: 2/3, 4-32
- A.24: 6 digit numbers, 1-53
- A.24: Add comment: these are tools, not recipes, 1-42
- A.24: An Introduction to R Sec. 10.8, A-3
- A.24: Check revise Tukey-Anscombe, 2-19
- A.24: Comment on Tukey vs Scheffé, 2-29
- A.24: Comment on trellis, 4-23
- A.24: Does memory help to increase the bound?, 4-15
- A.24: LB: meaning?, 1-24
- A.24: LB: well, all this depends on the number of data and the number of simulations., 1-25
- A.24: Melbourne data public?, 4-35
- A.24: Monte Carlo Coverage: work out as example, 1-26
- A.24: RB define codimension, 4-18
- A.24: RB preferably small: why?, 4-29
- A.24: RB: If you believe the distribution assumption, for some good reason, 3-15
- A.24: RB:  $\hat{\beta}$  may not be the only LSE or be unbiased, 2-6
- A.24: RB: is X now random?, 2-21
- A.24: Ref to curse of dimension, 4-24
- A.24: Remark on small setosa, 4-27
- A.24: Scheffe: fix lower/upper plotting bound, 2-29
- A.24: Simulation, Power, 3-21
- A.24: Stein shrinkage, 4-24
- A.24: Variance bias dilemma, 4-24
- A.24: add color, 4-27
- A.24: add legend, control colour, 4-17
- A.24: add magnitude, 4-24
- A.24: add minimal guaranteed effect, 3-31
- A.24: add more regression graphics, 2-1
- A.24: add power function, 3-15
- A.24: add power simulation - see Rnw.tex source, 3-16
- A.24: add tables, 4-42
- A.24: adjust scale, orientation, 4-32
- A.24: adjust style, 4-38
- A.24: better RExercise, 1-44
- A.24: ch.4: major rewrite, 4-1
- A.24: ch01: add examples f. outer, 1-25
- A.24: ch02: Add example, 2-41
- A.24: ch03: RExercise: power, 3-10
- A.24: ch03: work out, 3-10
- A.24: check contrasts, 2-23
- A.24: check iris inverse ratio, 4-28
- A.24: check plot, 2-28
- A.24: check size of barchart plot, 1-22
- A.24: clarify: information <-> assumptions, 3-15
- A.24: clean up discussion of factors, 2-23
- A.24: color R version, 4-15
- A.24: colour for left not visible in b/w, 3-5
- A.24: comment differences, 1-39
- A.24: conditional, function parameters, promises, 2-42
- A.24: define studentised residuals , 2-19
- A.24: discuss choice of models - model dependent estimation, 2-23
- A.24: discuss scale dependency, 4-52

## Index-12

A.24: example in appendix, 1-50  
A.24: formal variance bias, 4-29  
A.24: formats: add frame. This may be long. longtable environment?, Private-10  
A.24: formats: avoid page break, Private-6  
A.24: formats: show text in tables ragged right. Add hyphenation, Private-6  
A.24: give proper comment to external data access, 4-38  
A.24: improve Scheffe, 2-27  
A.24: improve example for names, A-7  
A.24: introduce noncentral t in ch. 2, 3-22  
A.24: mention unbiased, 2-8  
A.24: more precise word than easier, 1-41  
A.24: one sided or two sided?, 3-22  
A.24: or change plot co-ordinates?, 4-21  
A.24: reduce body fat case study drastically, 4-43  
A.24: reduce preliminary information on fat, 4-43  
A.24: ref appendix: reshape functions A:datamamip 2-24  
A.24: refer to reShape data frames in Hmisc Harrell 4.2.7, A-14  
A.24: remove line caps - bar line is too long: "h" for histogram-like vertical lines., 1-22  
A.24: rows, columns, lanes?, 2-35  
A.24: see comment LB, 1-32  
A.24: std error, t-test, point-wise confidence, 2-19  
A.24: studentised residuals, pointwise confidence intervals, 2-19  
A.24: supply better support for legend, colour key, 4-12  
A.24: this is first theorem. make numeration more logical, 1-22  
A.24: use TeX image placement, 1-53  
A.24: use attach, 4-29  
A.24: use principal components in modelling. pc regression?, 4-52  
A.24: www: add references, 1-40  
A.24: zum einen..zum anderen, 2-19

**aa2\_itd (alias2)**, Private-1  
added variable plots, **4-30**  
alpha channel, **1-53**

## Index

analysis of variance, **2-20**  
annotation, A-32  
argument  
    default, 1-3  
  
bandwidth, **1-11**  
Bonferroni correction, **2-10**  
bootstrap, **3-10**  
Box-Cox transformation, **2-36**  
brushing, **4-7**  
  
captions, A-32  
coefficient of variation, **3-27**  
conditioned, **4-18**  
contrast, **2-23, 2-30**  
curse of dimensionality, **4-38**  
  
data features, **1-33**, 3-30  
data structures, 1-19, A-9  
data types, 2-42  
date  
    see DateTimeClasses, A-26  
DateTimeClasses, A-26  
Debugging, 1-51, A-19  
design matrix, **2-3, 2-16**  
distribution  
    marginal, **4-8**  
  
entry, **Private-10**  
exact test, 3-15  
expected value, **1-35**  
  
factor, **2-5**  
    levels, 2-5  
fit, **2-2**  
function, **1-47**, A-17–A-21  
    polymorphic, *see polymorphic*  
  
Gauss-Markov estimator, **2-6**  
  
hat matrix, **2-8**  
histogram, **1-11**  
  
interquartile range, **1-36**  
  
join, *see merge*A-14

## Index

kernel, **1-11**  
least squares estimator, **2-6**  
legend, A-32  
linking, **4-7**  
**locfit**, *4-7*  
**loglin**, *1-31*  
  
**MASS**, *2-44, 2-45*  
missing  
    argument, *1-49, A-17*  
    number, *1-1*  
misssing  
    number, *2-43*  
model  
    linear, *2-3*  
    simple linear, **2-8**  
model function, **2-2**  
multidimensional scaling, **4-2**  
**mva**, *4-30*  
  
one sample case, **3-9**  
  
plot, *1-6*  
plotmath, *1-15*  
polymorphic, *1-4, 1-55, 2-43, 2-44, A-6*  
power, **1-34**  
principal component analysis, **4-30**  
probability plot, **1-41**  
Profiling, A-19  
projection pursuit, **4-14**  
  
**quantile**, **1-36**  
**quantile plot**, **1-41**  
  
random numbers, *1-4*  
    pseudo, *1-8*  
    reproducible, *A-37*  
random seed, *see* seed  
regression  
    linear, *2-3*  
regressor, **2-2**  
residual, **2-2, 2-8**  
response, **2-2**  
  
samples  
    repeated, **1-33**

## Index-13

scale  
    categorical, *2-5*  
    ordinal, *2-5*  
scale shift family, **3-5**  
scatterplot matrix, **4-8**  
search path, **1-55, 1-56, A-1**  
seed, **1-8, A-36**  
serial plot, **1-6**  
shift family, **3-4**  
standard deviation, **1-36**  
stochastically smaller, **3-4**  
substitute, **Private-10**  
  
test  
     $\chi^2$ , *1-31*  
    exact, *3-15*  
    Kolmogorov-Smirnov, *1-29*  
    median, *1-30*  
    midrange, *1-30*  
    Monte Carlo, **1-25**  
    t, *3-12, 3-15, 3-16, 3-21, 3-22, 3-25*  
    Wilcoxon, *3-15*  
tie, *3-15*  
time  
    *see* DateTimeClasses, A-26  
two sample case, **3-9**  
  
used, *Private-11*  
  
variance, **1-36**  
    empirical, **1-36**  
    residual, **2-9**  
    sample, **1-36**

Wilkinson-Rogers notation, **2-4, 2-17–2-18, 2-23, A-29–A-30**



---

# Index

---

From: /Users/sw/R/R-devel/doc/KEYWORDS

GROUPED Keywords

Group **Graphics**

Topic **aplot**

Add to Existing Plot / Internal Plot

**abline**, 1-13

**axis**, A-31

**contour**, 4-1, 4-2, A-30

Topic **dplot**

Computations Related to Plotting

Topic **hplot**

High-Level Plots

Topic **iplot**

Interacting with Plots

Topic **color**

Color, Palettes etc

Topic **dynamic**

Dynamic Graphics

Topic **device**

Graphical Devices

Group **Basics**

Topic **sysdata**

Basic System Variables

Topic **datasets**

Datasets Available by Data(.)

Topic **data**

Environments, Scoping, Packages

Topic **manip**

Data Manipulation

Topic **attribute**

Data Attributes

Topic **classes**

Data Types (not OO)

Topic **classes: character**

Character Data ("String") Operations

Topic **classes: complex**

Complex Numbers

Topic **classes: category**

Categorical Data

Topic **classes: NA**

Missing Values

Topic **list**

Lists

Topic **chron**

Dates and Times

Topic **package**

Package Summaries

Private-2

Index

## Group **Mathematics**

### Topic **array**

Matrices and Arrays

### Topic **array: algebra**

Linear Algebra

### Topic **arith**

Basic Arithmetic and Sorting

### Topic **math**

Mathematical Calculus etc.

### Topic **logic**

Logical Operators

### Topic **optimize**

Optimization

### Topic **symbolmath**

“Symbolic Math”, as Polynomials, Fractions

### Topic **graphs**

Graphs, (not Graphics), e.g. Dendograms

## Group **Programming, Input/Ouput, and Miscellaneous**

### Topic **programming**

Programming

### Topic **programming: interface**

Interfaces to Other Languages

### Topic **IO**

Input/output

### Topic **IO: file**

Files

### Topic **IO: connection**

Connections

### Topic **IO: database**

Interfaces to Databases

### Topic **iteration**

Looping and Iteration

### Topic **methods**

Methods and Generic Functions

### Topic **print**

Printing

### Topic **error**

Error Handling

### Topic **environment**

Session Environment

### Topic **internal**

Internal Objects (not Part of API)

### Topic **utilities**

Utilities

### Topic **misc**

Miscellaneous

### Topic **documentation**

Documentation

### Topic **debugging**

Debugging Tools

## Group **Statistics**

Index	Private-3
Topic <b>datagen</b> Functions for Generating Data Sets	Topic <b>smooth: loess</b> Loess Objects
Topic <b>distribution</b> Probability Distributions and Random Numbers	Topic <b>cluster</b> Clustering
Topic <b>univar</b> Simple Univariate Statistics	Topic <b>tree</b> Regression and Classification Trees
Topic <b>htest</b> Statistical Inference	Topic <b>survey</b> Complex Survey Samples
Topic <b>models</b> Statistical Models	Topic <b>classif</b> Classification
Topic <b>models: regression</b> Regression	Topic <b>spatial</b> Spatial Statistics
Topic <b>models: nonlinear</b> Non-linear Regression (only?)	Topic <b>neural</b> Neural Networks
Topic <b>robust</b> Robust/Resistant Techniques	
Topic <b>design</b> Designed Experiments	
Topic <b>multivariate</b> Multivariate Techniques	
Topic <b>ts</b> Time Series	
Topic <b>survival</b> Survival Analysis	
Topic <b>nonparametric</b> Nonparametric Statistics [w/o ‘Smooth’]	
Topic <b>smooth</b> Curve (and Surface) Smoothing	



---

## CHAPTER B

---

S formats. Remove for public version

---

### B.1 Spellings and Hyphenations

*t*-test reference: ?t.test

barchart lattice statistic: still to check usage. When to use “statistics”?

### B.2 Formats

**Example B.1** Nur *expl*

**Example B.1**: dem Schunk  
Vor dem Schunk

\$chunck:fo1

expl:fo1 expl:fo1 expl:fo2 Schunck:fo1  
B.I. expl:fo2 B.I. Schunk:fo1 B.I.

### B.3 Test R function index

*aa1* *itfun()*

*aa2* *itd*

---

T<sub>E</sub>X Example Start: exa2

T<sub>E</sub>X Example Input: exa2

\ixr{fun1}

T<sub>E</sub>X Example Output: exa2

exa2

Private-1

Private-2

TeX Example End: exa2

---

---

TeX Example Start: exa3

TeX Example Input: exa3

\ixr[fun2]{fun3}

TeX Example Output: exa3

exa3

fun2]fun3

TeX Example End: exa3

---

#### B.4 R styles

Sinput

---

TeX Example Start: exa4

TeX Example Input: exa4

```
Before private note
\begin{Rprivnote}
A private note
\end{Rprivnote}
After private note
```

TeX Example Output: exa4

Before private note

*Private note:*

*A private note*

exa4

After private note

TeX Example End: exa4

---

TeX Example Start: exa5

TeX Example Input: exa5

```
\begin{Sinput}
>plot(x) # "1+1"
\end{Sinput}
```

TeX Example Output: exa5

exa5

*Input*

TeX Example End: exa5

---

Soutput

---

TeX Example Start: exa6

TeX Example Input: exa6

```
\begin{Soutput}
abc <- "1+1"
\end{Soutput}
```

TeX Example Output: exa6

exa6

*Output*

TeX Example End: exa6

---

Sinput and Soutput in an Schunk

---

TeX Example Start: exa7

TeX Example Input: exa7

Private-4

```
\begin{Schunk}
\begin{Sinput}
>plot(x) # "1+1"
\end{Sinput}
\begin{Soutput}
abc <- "1+1"
\end{Soutput}
\end{Schunk}
```

TEX Example Output: exa7

plot(x) # "1+1"

exa7

*Input*

*Output*

TEX Example End: exa7

irin

TEX Example Start: exa8

TEX Example Input: exa8

```
\begin{irin}
abc$d <- "1+1" #R input code block
\end{irin}
```

TEX Example Output: exa8

exa8

abc\$d <- "1+1" #R *input code block*

TEX Example End: exa8

irout

---

T<sub>E</sub>X Example Start: exa9

T<sub>E</sub>X Example Input: exa9

```
\begin{irout}
abc <- "1+1" #R output code block
\end{irout}
```

T<sub>E</sub>X Example Output: exa9

---

exa9      abc <- "1+1" #R output code block

---

T<sub>E</sub>X Example End: exa9

---

---

Scode

---

T<sub>E</sub>X Example Start: exa10

T<sub>E</sub>X Example Input: exa10

```
\begin{Scode}
abc <- "1+1" # R output code block in Scode environment
\end{Scode}
```

T<sub>E</sub>X Example Output: exa10

---

exa10      abc <- "1+1" # R output code block in Scode environment

---

T<sub>E</sub>X Example End: exa10

---

---

T<sub>E</sub>X Example Start: exa11

T<sub>E</sub>X Example Input: exa11

A sample \ircode{abc <- "1+1"} inline code

## Private-6

TEX Example Output: exa11

exa11  
A sample `abc <- "1+1"` inline code

### B.5 Tables

#### ToDo:

formats:  
show text  
in tables  
ragged  
right. Add  
hyphen-  
ation

#### ToDo:

formats:  
avoid page  
break

TEX Example Start: exa12

TEX Example Input: exa12

```
\begin{tabthree}{\R{} Sequenzen}{}{}\n\\txdesc{\\irfunx{manip}{c}}{''combine''. \\ Kombiniert Argumente zu einem neuen Vektor.}%\n\\txsynt{\\ircode{c(..., recursive=FALSE)}}\n\\txexpls{\\ircode{c(1,2,3)}\\newline\\ircode{c(x,y)}}\n%\n\\txdesc{\\ircode{:}}{Erzeugt ganzzahlige Sequenz von \\ircode{Anfang} bis \\ircode{Ende}.}\n\\txsynt{\\ircode{Anfang:Ende}}\n\\txexpls{\\ircode{1:10}\\newline\\ircode{10:1}}\n%\n\\txdesc{\\irfunx{manip}{seq}}{Erzeugt allgemeine Sequenzen.}\n\\txsynt{Siehe \\ircode{help(seq)}}\n%\n\\txdesc{\\irfunx{manip}{rep}\\ keyword{chron}{rep}}{Wiederholt Argument.}\n\\txsynt{\\ircode{rep(x, times, ...)}}\n\\txexpls{\\ircode{rep(x,3)}\\newline\\ircode{rep(1:3,c(2,3,1))}}\n\\end{tabthree}
```

TEX Example Output: exa12

R Sequenzen	
<code>c()</code>	"combine". Kombiniert Argumente zu einem neuen Vektor. Syntax: <code>c(..., recursive=FALSE)</code> Examples: <code>c(1,2,3)</code> <code>c(x,y)</code>

(cont.)→

<b>R Sequenzen</b> (cont.)	
:	Erzeugt ganzzahlige Sequenz von <i>Anfang</i> bis <i>Ende</i> . Syntax: <i>Anfang:Ende</i> Examples: <i>1:10</i> <i>10:1</i>
<i>seq()</i>	Erzeugt allgemeine Sequenzen. Syntax: Siehe <i>help(seq)</i>
<i>rep()</i>	Wiederholt Argument. Syntax: <i>rep(x, times, ...)</i> Examples: <i>rep(x,3)</i> <i>rep(1:3, c(2,3,1))</i>

**exa12**

---

TeX Example End: exa12

---

---

---

TeX Example Start: exa13**TeX Example Input: exa13**

```
\begin{RExercise}[Optional Title]{1}
&abc\\
de&fg\\
\hline
\end{RExercise}
```

---

**TeX Example Output: exa13**

---

**1**

Exercise B.1	Optional Title
	abc
de	fg

**exa13**

---

**TeX Example End: exa13**

---

## Private-8

---

TeX Example Start: exa14

TeX Example Input: exa14

```
\begin{tabthree}{Verteilungen}{Kurzname}{Defaults}
Beta & beta & \ircode{shape1, \, shape2, \, ncp=0} \\ \hline
Binomial & binom & \ircode{size, prob} \\ \hline
Cauchy & cauchy & \ircode{location = 0, scale = 1} \\ \hline
\chi^2 & chisq & \ircode{df, ncp=0} \\ \hline
Exponential & exp & \ircode{rate = 1} \\ \hline
F & f & \ircode{df1, df2, ncp = 0} \\ \hline
Gamma & gamma & \ircode{shape, scale = 1} \\ \hline
Gauss & norm & \ircode{mean = 0, sd = 1} \\ \hline
Geometrisch & geom & \ircode{prob} \\ \hline
Hypergeometrisch & hyper & \ircode{m, n, k} \\ \hline
Lognormal & lnorm & \ircode{meanlog = 0, sdlog = 1} \\ \hline
Logistisch & logis & \ircode{location = 0, scale = 1} \\ \hline
Negativ-Binomial & nbinom & \ircode{size, prob} \\ \hline
Poisson & pois & \ircode{lambda} \\ \hline
Student's t & t & \ircode{df} \\ \hline
Tukey Studentised Range & tukey & \\ \hline
Uniform & unif & \ircode{min = 0, max = 1} \\ \hline
Wilcoxon Signed Rank & signrank & \ircode{n} \\ \hline
Wilcoxon Rank Sum & wilcox & \ircode{m, n} \\ \hline
Weibull & weibull & \ircode{shape, scale = 1}
\\
\end{tabthree}
```

TeX Example Output: exa14

<i>Verteilungen</i>	<i>Kurzname</i>	<i>Defaults</i>
Beta	beta	<i>shape1, shape2, ncp=0</i>
Binomial	binom	<i>size, prob</i>
Cauchy	cauchy	<i>location = 0, scale = 1</i>
$\chi^2$	chisq	<i>df, ncp=0</i>
Exponential	exp	<i>rate = 1</i>
$F$	$f$	<i>df1, df2, ncp = 0</i>
Gamma	gamma	<i>shape, scale = 1</i>
Gauss	norm	<i>mean = 0, sd = 1</i>
Geometrisch	geom	<i>prob</i>
Hypergeometrisch	hyper	<i>m, n, k</i>

(cont.)→

<i>Verteilungen</i>	<i>Kurzname</i>	<i>Defaults</i>
Lognormal	lnorm	<i>meanlog = 0, sdlog = 1</i>
Logistisch	logis	<i>location = 0, scale = 1</i>
Negativ-Binomial	nbinom	<i>size, prob</i>
Poisson	pois	<i>lambda</i>
Student's t	t	<i>df</i>
Tukey Studentised Range	tukey	
Uniform	unif	<i>min = 0, max = 1</i>
Wilcoxon Signed Rank	signrank	<i>n</i>
Wilcoxon Rank Sum	wilcox	<i>m, n</i>
Weibull	weibull	<i>shape, scale = 1</i>

exa14

---

TeX Example End: exa14

---

## B.6 Help

Help information is copied as is from the R help files. **Do not change format or spelling in the help information.**

---

TeX Example Start: exa15

---

TeX Example Input: exa15

```
%\noindent\begin{longtable}{|p{.98\linewidth}|}
\begin{RHelp}
\begin{RHelpInput}
help(lm)
\end{RHelpInput}
\todo[formats: add frame. This may be long. longtable environment?]
\begin{RHelpOutput}{help(lm)}{lm} \hfill package:base \hfill
... see help test file
\end{RHelpOutput}
\end{RHelp}
%\end{longtable}
```

R Documentation\\Fitting

Private-10

---

TEX Example Output: exa15

---

**ToDo:**  
formats:  
add frame.  
This may  
be **exa15**,  
longtable  
environ-  
ment?

---

help(lm) help(lm)lm package:base R Documentation  
Fitting Linear Models

---

... see help test file

---

TEX Example End: exa15

---

## B.7 Indexed Terms

---

TEX Example Start: exa16

---

TEX Example Input: exa16

---

An index \itd[{substitute}]{entry} definition with substitute.\\"

---

TEX Example Output: exa16

---

An index *entry* definition with substitute.

---

TEX Example End: exa16

---

---

TEX Example Start: exa17

---

---

TEX Example Input: exa17

---

An index \itd{entry} definition without substitute.\\"

---

TEX Example Output: exa17

---

An index *entry* definition without substitute.

**exa17**TeX Example End: exa17

---

---

TeX Example Start: exa18

**exa18** **Tex Example Input: exa18**

A \itx{used} index entry\\

**Tex Example Output: exa18**A *used* index entry**exa18** **Tex Example End: exa18**

---

---



---

## CHAPTER C

# S help — remove for public version

---

### C.1 Ch 01

---

---

**help(runif)**

---

**Uniform**

*The Uniform Distribution*

---

#### Description

These functions provide information about the uniform distribution on the interval from `min` to `max`. `dunif` gives the density, `punif` gives the distribution function `qunif` gives the quantile function and `runif` generates random deviates.

#### Usage

```
dunif(x, min=0, max=1, log = FALSE)
punif(q, min=0, max=1, lower.tail = TRUE, log.p = FALSE)
qunif(p, min=0, max=1, lower.tail = TRUE, log.p = FALSE)
runif(n, min=0, max=1)
```

#### Arguments

<code>x,q</code>	vector of quantiles.
<code>p</code>	vector of probabilities.
<code>n</code>	number of observations. If <code>length(n) &gt; 1</code> , the length is taken to be the number required.
<code>min,max</code>	lower and upper limits of the distribution. Must be finite.
<code>log, log.p</code>	logical; if TRUE, probabilities p are given as log(p).
<code>lower.tail</code>	logical; if TRUE (default), probabilities are $P[X \leq x]$ , otherwise, $P[X > x]$ .

### Details

If `min` or `max` are not specified they assume the default values of `0` and `1` respectively.

The uniform distribution has density

$$f(x) = \frac{1}{max - min}$$

for  $min \leq x \leq max$ .

For the case of  $u := min == max$ , the limit case of  $X \equiv u$  is assumed, although there is no density in that case and `dunif` will return `NaN` (the error condition).

`runit` will not generate either of the extreme values unless `max = min` or `max-min` is small compared to `min`, and in particular not for the default arguments.

### References

Becker, R. A., Chambers, J. M. and Wilks, A. R. (1988) *The New S Language*. Wadsworth & Brooks/Cole.

### See Also

`.Random.seed` about random number generation, `rnorm`, etc for other distributions.

### Examples

```
u <- runif(20)

## The following relations always hold :
punif(u) == u
dunif(u) == 1

var(runit(10000))#- ~ = 1/12 = .08333
```

---

---

## CHAPTER D

# To Do

---

### D.1 For release

Fix colour error at help page(?)  
Fix runoff text in too long input lines: ex:1.3  
Check/redo image sizes  
Fix " ->  $\beta$  in help(loess) family

### D.2 misc

add function keywords using /Users/sw/R/R-devel/doc/KEYWORDS  
discuss name extension  
in Sframe: remove vertical space after Beispiel.  
summar f. typeof  
add margin.table, addmargins Verzani p. 72-73  
add stripchart Verzani p. 79 add jitter Verzani p.90  
Check ' ToDo' entries in index!!!  
Fix refs at exercises  
Fine tune graphics size  
Add bibliography by chapter  
Center email/URL on title page  
Grid/Lattice Check lattice vs par usage  
Classes?  
Check libraries/packages  
Fix Bibliography to german header & special words  
Check missing fonts: grep 'LaTeX Font Info' \*.log

Private-2

TO DO

```
$Source: /u/math/j40/cvsroot/lectures/src/SIntro/Rintro/chapters/Rintro_main.tex,v $
$Revision: 1.16 $
$Date: 2008/06/17 07:24:15 $
$Name:  $
$Author: j40 $

$Source: /u/math/j40/cvsroot/lectures/src/SIntro/Rintro/Rintro_temp.tex,v $
$Revision: 1.13 $
$Date: 2008/06/18 08:20:05 $
$Name:  $
$Author: j40 $
```