

Computational Statistics

An Introduction to 

Günther Sawitzki

 CRC Press
Taylor & Francis Group

A CHAPMAN & HALL BOOK

Computational Statistics

An Introduction to Exercises

April 30, 2014 DRAFT

Günther Sawitzki
StatLab Heidelberg
April 30, 2014

in preparation

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

Contents

1	Basic Data Analysis	1
1.1	R Programming Conventions	1
1.2	Generation of Random Numbers and Patterns	1
1.2.1	Random Numbers	1
1.2.2	Patterns	1
1.3	Case Study: Distribution Diagnostics	2
1.3.1	Distribution Functions	2
1.3.2	Histograms	2
	Barcharts	2
1.3.3	Statistics of Distribution Functions; Kolmogorov-Smirnov Tests	2
	Monte Carlo Confidence Bands	3
1.3.4	Statistics of Histograms and Related Plots; χ^2 -Tests	6
1.4	Moments and Quantiles	8
1.5	R Complements	11
1.5.1	Random Numbers	11
1.5.2	Graphical Comparisons	11
1.5.3	Functions	12
	Vectorisation	13
	Compilation	13
1.5.4	Enhancing Graphical Displays	13
1.5.5	R Internals	13
	Executing Files	13
1.5.6	Search Paths, Frames and Environments	13
1.6	Additional Exercises	13

1.6.1	Packages	15
	Using Packages	15
	Building Packages	15
	Compilation	15
1.7	Statistical Summary	15
1.8	Literature and Additional References	15
2	Regression	17
2.1	General Regression Model	17
2.2	Linear Model	17
2.2.1	Factors	17
2.2.2	Least Squares Estimation	17
2.2.3	Regression Diagnostics	18
2.2.4	More Examples for Linear Models	19
2.2.5	Model Formulae	19
2.2.6	Gauss-Markov Estimator	19
2.3	Variance Decomposition and Analysis of Variance	20
2.4	Simultaneous Inference	23
2.4.1	Scheffé's Confidence Bands	23
2.4.2	Tukey's Confidence Intervals	23
	Case Study: Titre Plates	23
2.5	Beyond Linear Regression	23
	Transformations	23
2.5.1	Generalised Linear Models	23
2.5.2	Local Regression	24
2.6	R Complements	24
2.6.1	Discretisation	24
2.6.2	External Data	24
2.6.3	Testing Software	24
2.6.4	R Data Types	25
2.6.5	Classes and Polymorphic Functions	25
2.6.6	Extractor Functions	25
2.7	Statistical Summary	25
2.8	Literature and Additional References	25

3	Comparisons	27
3.1	Shift/Scale Families, and Stochastic Order	29
3.2	<i>QQ</i> Plot, <i>PP</i> Plot, and Comparison of Distributions	31
3.2.1	Kolmogorov-Smirnov Tests	33
3.3	Tests for Shift Alternatives	34
3.4	A Road Map	35
3.5	Power and Confidence	35
3.5.1	Theoretical Power and Confidence	35
3.5.2	Simulated Power and Confidence	35
3.5.3	Quantile Estimation	35
3.6	Qualitative Features of Distributions	36
3.7	Statistical Summary	36
3.8	Literature and Additional References	36
4	Dimensions 1, 2, 3, . . . , ∞	37
4.1	R Complements	37
4.2	Dimensions	37
4.3	Selections	37
4.4	Projections	37
4.4.1	Marginal Distributions and Scatter Plot Matrices	37
4.4.2	Projection Pursuit	37
4.4.3	Projections for Dimensions 1, 2, 3, . . . 7	39
4.4.4	Parallel Coordinates	39
4.5	Sections, Conditional Distributions and Coplots	39
4.6	Transformations and Dimension Reduction	45
4.7	Higher Dimensions	45
4.7.1	Linear Case	45
	Partial Residuals and Added Variable Plots	45
4.7.2	Non-Linear Case	45
	Example: Cusp Non-Linearity	45
4.7.3	Case Study: Melbourne Temperature Data	45
4.7.4	Curse of Dimensionality	45
4.7.5	Case Study: Body Fat	45
4.8	High Dimensions	48
4.9	Statistical Summary	48

viii

References	49
Functions and Variables by Topic	51
Function and Variable Index	51
Subject Index	51

CHAPTER 1

Basic Data Analysis

1.1 R Programming Conventions

1.2 Generation of Random Numbers and Patterns

1.2.1 Random Numbers

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 $(0, 1)$?</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 clearly as possible.</p> <p><i>Hint:</i> For comparison, you can keep several plots in a window. The code</p> <pre data-bbox="724 1487 1002 1514">par(mfrow = c(2, 3))</pre> <p>parametrises the graphics system to show six plots simultaneously, arranged row wise as a 2×3 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>

1.2.2 Patterns

Exercise 1.2	
	<p>Use</p> <pre style="text-align: center;"><code>plot(sin(1:100))</code></pre> <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>

1.3 Case Study: Distribution Diagnostics

1.3.1 First Pass for Example ?? : Distribution Functions

1.3.2 First Pass for Example ?? : Histograms

Exercise 1.3	
	<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 $(0, 0.5]$, $(0.5, 1)$.</p> <pre style="text-align: center;"><code>hist(runif(100), breaks = c(0, 0.5, 1))</code></pre> <p>Repeat the experiment with random numbers generated by <code>rnorm(100)</code> and compare the results from <code>runif(100)</code> and <code>rnorm(100)</code>.</p>

Exercise 1.4	
	<p>Modify Example ?? (page ??) to include the kernel name and the bandwidth used in the kernel density estimation.</p> <p>You have to store the result from <code>density()</code> and access its components in analogy to Example ?? (page ??).</p>

Barcharts

1.3.3 Statistics of Distribution Functions; Kolmogorov-Smirnov Tests

Exercise 1.5	
	Using <code>help(rbeta)</code> you get information about the functions available to work with beta distributions. Generate plots for the densities of the beta distribution for $n = 16, 32, 64, 128$ and $i = n/4, n/2, 3n/4$. Use the function <code>curve()</code> to generate the plots. For more information, see <code>help(curve)</code> .

Exercise 1.6	
	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 $X_{(i)}$ an adequate reference? Are there alternatives that can serve as references? If you see alternatives, give an implementation.

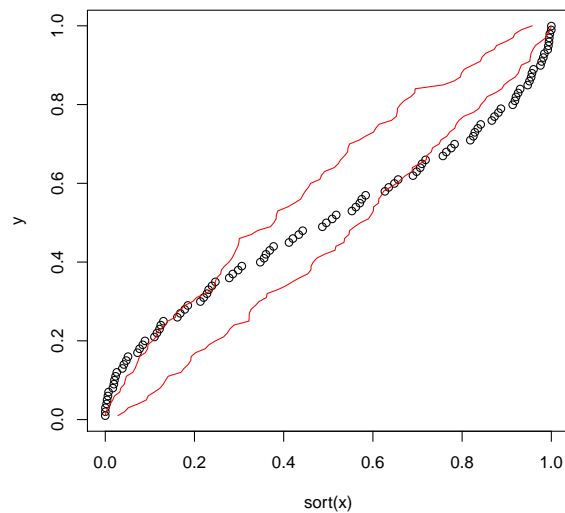
Monte Carlo Confidence Bands

Example 1.1: Monte Carlo Confidence Bands

```

Input
x <- (sin(1:100)+1)/2 # demo example only
y <- (1:length(x))/length(x)
plot(sort(x), y)
nrsamples <- 19 # no of simulations
samples <- matrix(data = 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")

```



Exercise 1.7	
	Make use of the <code>help()</code> -function and comment on Example 1.1 step by step. Take special note of the new functions that are introduced here.

R Iterators	
--------------------	--

(cont.)→

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

Exercise 1.8	
*	<p>Why 19?</p> <p><i>Hint:</i> Try to take an abstract simplified view of the problem first: let T be a measurable function and $X_0, X_1, \dots, X_{nrsamples}$ independent samples with a common distribution function.</p> <p>What is $P(T(X_0) > T(X_i))$ for all $i > 0$?</p> <p>In a second step, give an abstract formulation for the example above. Then consider the special case $nrsamples = 19$.</p>

Exercise 1.9	Monte Carlo Coverage
*	<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 sim 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 sim 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>

Theorem 1.1 For all integer n and any positive λ , we have

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

Proof. [6], Corollary 1 \square

This inequality is valid even if F is not continuous.

Exercise 1.10	Finite Sample Bounds
	Use the inequality given in Theorem 1.1 to calculate bounds for $\sqrt{n} \sup F_n - F $.
	Add finite sample bands to the empirical distribution function.

Exercise 1.11	
	Using <code>help(ks.test)</code> you get information on how to invoke the function <code>ks.test</code> . Which results do you expect if you test the following vectors for a uniform distribution? <div style="text-align: center;"> <code>1:100</code> <code>runif(100)</code> <code>sin(1:100)</code> <code>rnorm(100)</code> </div> Perform these tests and discuss the results. For the test, scale the values so that they fall into the interval $[0, 1]$, or use a uniform distribution on an interval that is adapted to the data.

1.3.4 Statistics of Histograms and Related Plots; χ^2 -Tests

Exercise 1.12	
	Use <code>help(chisq.test)</code> to see the calling structure for χ^2 tests. Apply it to test the hypothesis ($p_j = 1/J$), $J = 5$ on the following vectors of bin counts: <div style="text-align: center;"> <code>(3 3 3 3 3)</code> <code>(1 2 5 3 3)</code> <code>(0 0 9 0 6)</code>. </div>

Exercise 1.13	
	Which results do you expect if you use a χ^2 test to check the following vectors for a uniform distribution? <div style="text-align: center;"> <code>1:100</code> <code>runif(100)</code> </div> <div style="text-align: right;">(cont.)→</div>

Exercise 1.13	(cont.)
	<pre data-bbox="754 369 895 439">sin(1:100) rnorm(100)</pre> <p data-bbox="480 454 983 483">Perform these tests and discuss the results.</p> <p data-bbox="480 499 1254 618"><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 component of its result.</p>

Exercise 1.14	
*	<p data-bbox="480 757 1254 819">Sketch comparable test environments for fixed and adaptive choice of histogram cells.</p> <p data-bbox="480 835 1254 927">For fixed and for adaptive choice of histogram cells draw $s = 1000$ samples of size 50 from <code>runif()</code>. Calculate in both settings the formal χ^2 statistics and plot its distribution functions.</p> <p data-bbox="480 943 895 972">Compare the distribution functions.</p>

Exercise 1.15	
	<p data-bbox="480 1106 1254 1169">For $n = 10, 50, 100$, draw 300 samples using <code>runif(n)</code>. For each sample, calculate the χ^2 and the Kolmogorov-Smirnov statistic.</p> <p data-bbox="480 1184 1078 1214">You have to choose a χ^2 test. What is your choice?</p> <p data-bbox="480 1229 1254 1292">Plot the distribution functions of these statistics and compare them to the theoretical (asymptotic) distributions.</p> <p data-bbox="480 1308 1254 1370">Are there any indications against the assumption of independent uniform random numbers?</p> <p data-bbox="480 1386 1254 1471"><i>Hint:</i> The functions for the χ^2 and Kolmogorov-Smirnov test keep their internal information as a list. To get the names of the list elements, you can create a sample object. For example, use</p> <pre data-bbox="659 1487 1070 1516">names(chisq.test(runif(100))).</pre>

Exercise 1.16	
**	<p>Analyse the power of the Kolmogorov-Smirnov test and the χ^2 tests. Select values for n, m and α, and choose 9 pairs for (a, b). 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 χ^2 test with 10 cells of equal size on $(0, 1)$.</p>
	<p>Choose alternative parameters (a, b) so that you can compare the decision rules along the following lines:</p> <ul style="list-style-type: none"> i) $a = b$ ii) $b = 1$ iii) $a = 1$ <p>and run these simulations.</p>
	<p>Choose alternative parameters (a, b) so that you can compare the decision rules over the range $0 < a, b < 5$.</p> <p>Your conclusions?</p> <p><i>Hint: <code>outer(x, y, fun)</code> applies a function <code>fun()</code> to all pairs of values from x, y and returns the result as a matrix.</i></p> <p>Using</p> <p style="text-align: center;"><code>contour()</code></p> <p>you can generate a contour plot. See <code>demo("graphic")</code>.</p>

Exercise 1.17	
**	<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"> i) $x = 1..100 \bmod m$ for convenient m ii) $\sin(x)$ $x = 1..100$ iii) ... <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

Exercise 1.18	
	<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>

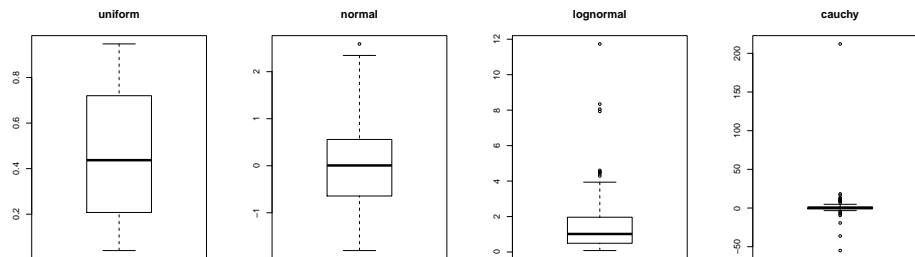
Exercise 1.19	
	<p>Generate a sample of 100 random variables from the distributions of Exercise 1.18.</p> <p>Estimate the median, and the lower and upper quartiles.</p> <p>Repeat the estimation for 1000 samples. Analyse the distribution of the estimated median, lower and upper quartiles from repeated samples.</p>

Example 1.2: 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)

```



Exercise 1.20	
	Modify Example 1.2 so that the plots are comparable: adjust the location so that the medians are at the same height. Adjust the scales so that the inter-quartile ranges have same length.

Exercise 1.21	
	<p>For continuous distributions and the median X_{med} we have</p> $P(X_i \geq X_{med}) = 0.5.$ <p>Hence we can find a k such that</p> $k = \min\{k : P(X_{(k)} \leq X_{med}) < \alpha\}$ <p>and $X_{(k)}$ as an upper bound for the median with confidence level $1 - \alpha$.</p> <p>Use this idea to construct a confidence interval for the median with confidence level $1 - \alpha = 0.9$.</p>
	<p>Modify the box-and-whisker plot to show this interval.</p> <p style="text-align: right;">(cont.)→</p>

Exercise 1.21	(cont.)
	<i>Hint:</i> You need the distribution function F_X , evaluated at the position marked by the order statistic $X_{(k)}$. The distributions of $F_X(X_{(k)})$ are discussed in Theorem ??.
	The box-and-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> .
*	Use an analogous strategy to get a distribution-independent confidence interval for the inter-quartile range.
***	Augment the box-and-whisker plot so that it gives information about the scale in a way that is statistically reliable. <i>Hint:</i> Why is it not sufficient to mark confidence intervals for the quartiles?

1.5 R Complements

1.5.1 Random Numbers

1.5.2 Graphical Comparisons

Exercise 1.22	
	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 of 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 ν be so that the t distribution is barely different from the normal distribution in these plots?
	How large must ν be so that the t distribution is barely different from the normal distribution if you compare the graphs of the distribution functions?

Exercise 1.23	
	Use <i>PP</i> plots instead of distribution functions to illustrate the χ^2 - and Kolmogorov-Smirnov approximations.

Exercise 1.24	
	Use <i>QQ</i> plots instead of distribution functions. Can you add confidence regions to these plots with the help of the χ^2 - resp. Kolmogorov-Smirnov statistics?

Exercise 1.25	
	Generate a matrix of dimensions $(nrow * ncol - 1)$, $length(x)$ with random numbers and use <code>apply()</code> to avoid the loop. <i>Hint:</i> See Example 1.1 (page 3).

Exercise 1.26	
	Use <code>rnorm()</code> to generate with 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.
	Add Monte Carlo bands from the envelope 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 coordinate 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 <code>qqnorm()</code> .
*	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

Exercise 1.27	
	Rework your programming exercises and write reusable parts as functions.

Exercise 1.28	
	<p>Write as functions:</p> <ul style="list-style-type: none"> • A function <code>ehist</code> showing an augmented histogram. • A function <code>eecdf</code> showing the empirical distribution. • A function <code>eqqnorm</code> showing a <i>QQ</i> plot with the standard normal distribution as comparison. • A function <code>eboxplot</code> showing a box-and-whisker plot. <p>and</p> <ul style="list-style-type: none"> • A wrapper function <code>eplot</code> showing a plot matrix with these four plots. <p>Your functions should call the standard functions (or modify them, if necessary) and guarantee that the plots have an adequate complete annotation.</p>

Vectorisation

I

Exercise 1.29	Vectorisation
	Write <code>sqrt0()</code> as a vectorised function using <code>ifelse()</code> .

Compilation

1.5.4 Complements: Enhancing Graphical Displays

Exercise 1.30	
	Use <code>help(plot)</code> to inspect the possibilities of customising the plot function. Information on details of the parameters is only available if you use <code>help(plot.default)</code> . Modify your latest plot so it has a correct main title.

1.5.5 Complements: R Internals

Executing Files

1.5.6 Search Paths, Frames and Environments

1.6 Additional Exercises

Exercise 1.31	Feature Detection								
	<p>This series of example tries to judge feature detection sensitivity of various displays for univariate data. for a preparation, find (and fix) a display arrangement that is convenient for your display, (for example <code>par(mfrow=c(5,4))</code>). Select a false detection rate you are willing to tolerate (for example, $2/20 = 10\%$).</p>								
*	<p>Write a function <code>plotdens <- function(n)</code> that draws n normal random numbers for each of the display frames.</p> <p>Find a number n_{sym} so that with $n \geq n_{sym}$ observations most results appear symmetric (i.e. the non-symmetric samples are below tolerance rate).</p> <p>Find a number $n_{unimodal}$ so that with $n \geq n_{unimodal}$ observations most results appear unimodal (i.e. the multimodal samples are below tolerance rate).</p>								
	<p>Modify your function by adding an additional parameter <code>plotdens <- function(n, generator=rnorm)</code> that allows to select a random number generator.</p> <p>For the following distributions, find a sample size that allow detection of the given features reliably within tolerance.</p> <table border="1" data-bbox="555 1010 1182 1178"> <thead> <tr> <th data-bbox="555 1010 699 1048">Generator</th> <th data-bbox="699 1010 1182 1048">Features</th> </tr> </thead> <tbody> <tr> <td data-bbox="555 1048 699 1086">uniform</td> <td data-bbox="699 1048 1182 1086">flat density</td> </tr> <tr> <td data-bbox="555 1086 699 1124">lognormal</td> <td data-bbox="699 1086 1182 1124">unimodal, skewed, long tail</td> </tr> <tr> <td data-bbox="555 1124 699 1178">Cauchy</td> <td data-bbox="699 1124 1182 1178">symmetric, frequent outliers, long tail.</td> </tr> </tbody> </table>	Generator	Features	uniform	flat density	lognormal	unimodal, skewed, long tail	Cauchy	symmetric, frequent outliers, long tail.
Generator	Features								
uniform	flat density								
lognormal	unimodal, skewed, long tail								
Cauchy	symmetric, frequent outliers, long tail.								

Exercise 1.32	Distribution
	<p>Prepare a plot with one display frame showing a test sample, the others showing uniform random samples. What is the required sample size to identify a normal distribution with an error rate below tolerance level? With this sample size, what is the false detection rate if you start with a uniform sample.</p>
	<p>Exchange the roles. Use a uniform sample as a test sample, and normal samples for comparisons. What is the required sample size to identify a uniform distribution with an error rate below tolerance level? With this sample size, what is the false detection rate if you start with a gaussian sample.</p>

Exercise 1.33	
	Repeat the example above using a QQ-Plot.
	Repeat the example above using a histogram.

1.6.1 Complements: Packages

Using Packages

Building Packages

Exercise 1.34	
	<p>Install the functions from Exercise 1.28 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 x, the statement <code>prompt(x)</code> generates a skeleton upon which you can build a documentation for x.</p>

Compilation

1.7 Statistical Summary

1.8 Literature and Additional References

[9] R Development Core Team (2000–2008): Writing R Extensions.

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

[12] Shorack, G. R.; Wellner, J. A.: *Empirical Processes with Applications to Statistics*. Wiley, New York, 1986.

[3] Gänßler, P.; Stute, W.: *Wahrscheinlichkeitstheorie*. Springer, Heidelberg, 1977.

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

CHAPTER 2

Regression

2.1 General Regression Model**2.2 Linear Model***2.2.1 Factors**2.2.2 Least Squares Estimation*

For this example data set, we get the least squares estimator using

Example 2.1: Least Squares Estimator	
<i>Input</i>	<code>lm(y ~ x)</code>
<i>Output</i>	Call: <code>lm(formula = y ~ x)</code> Coefficients: (Intercept) x -1.295 2.541

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.

Example 2.2: Linear Model Summary

```

----- Input -----
summary(lm( y ~ x))
----- Output -----
Call:
lm(formula = y ~ x)

Residuals:
    Min       1Q   Median       3Q      Max
-28.2790  -6.3290   0.8087   7.3920  28.3508

Coefficients:
              Estimate Std. Error t value Pr(>|t|)
(Intercept)  -1.2947     1.9719  -0.657   0.513
x             2.5410     0.0339  74.958 <2e-16 ***
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 9.785 on 98 degrees of freedom
Multiple R-squared:  0.9829,    Adjusted R-squared:  0.9827
F-statistic: 5619 on 1 and 98 DF,  p-value: < 2.2e-16

```

Exercise 2.2

Analyse the output of `lm()` shown in Example 2.2. 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.

*2.2.3 Regression Diagnostics***Exercise 2.3**

Let

$$yy \leftarrow 2.5 * x + 0.01 * x^2 + \text{err}$$

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?

Exercise 2.4	
	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?

2.2.4 More Examples for Linear Models

2.2.5 Model Formulae

Exercise 2.5	
	Write the four models from Section 2.2.4 using the R formula notation.
	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.

Exercise 2.6	
	<p>Generate three vectors of random variables with an $N(\mu_j, 1)$ distribution, $\mu_j = j$, $j = 1, 3, 9$, each of length 10, and combine these into a vector y.</p> <p>Generate a vector x with the values j, $j = 1, 3, 9$, each repeated 10 times.</p> <p>Calculate the Gauss-Markov estimator in the linear models $y \sim x$ and $y \sim \text{factor}(x)$.</p> <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>

2.2.6 Gauss-Markov Estimator

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

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.

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

2.3 Variance Decomposition and Analysis of Variance

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

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

The analysis of variance gives another representation and interpretation of linear models. For example, the regression result from (Example ??) gives the following analysis of variance representation:

<i>Example 2.3: Linear Model ANOVA Summary</i>							
				Input			
<code>summary(aov(lmres))</code>							
				Output			
	Df	Sum Sq	Mean Sq	F	value	Pr(>F)	
x	1	538022	538022	5619	<2e-16	***	
Residuals	98	9384	96				

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

`anova()` can be used to compute analysis of variance tables for one or more fitted model objects.

<i>Example 2.4: Linear Model ANOVA</i>							
				Input			
<code>anova(lm(y~x), lm(y~0+x))</code>							
				Output			
Analysis of Variance Table							
Model 1: y ~ x							
Model 2: y ~ 0 + x							
	Res.Df	RSS	Df	Sum of Sq	F	Pr(>F)	
1	98	9384.0					
2	99	9425.2	-1	-41.283	0.4311	0.513	

Exercise 2.12	
	Analyse the output of <code>lm()</code> shown in Example 2.2 (page 17). Which terms can you interpret now? Give a written report. For which terms do you need more information?

Exercise 2.13	One-Way Anova
*	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 2.14	Kiwi Hopp
	The industrial enterprise Kiwi Inc. ¹ 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 ² from a fixed starting height (ca. 2m). Figure 2.1, page 26, shows a design drawing. 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?
	Perform 30 test flights with a prototype and measure the time in 1/100s. (You will have to cooperate 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?
	Your task is to select a design for production. The variants under discussion are: <div style="margin-left: 40px;">rotor width 45mm</div> <div style="margin-left: 40px;">rotor width 35mm</div> <div style="margin-left: 40px;">rotor width 45mm with an additional fold for stabilisation</div> <div style="margin-left: 40px;">rotor width 35mm with an additional fold for stabilisation</div> <div style="text-align: right;">(cont.)→</div>

¹ Following an idea of Alan Lee, Univ. Auckland, New Zealand.

² Kiwis cannot fly.

Exercise 2.14	Kiwi Hopp (cont.)
	<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"> • a list of the observed data and a description of the experimental procedure • suitable plots of the data for each of the designs • an analysis of variance • a clear summary of your conclusions <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>
	<p>The fold will result in additional production cost. Give an estimate of the gain that can be achieved by this additional cost.</p>

2.4 Simultaneous Inference

2.4.1 Scheffé's Confidence Bands

2.4.2 Tukey's Confidence Intervals

Case Study: Titre Plates

2.5 Beyond Linear Regression

Transformations

2.5.1 Generalised Linear Models

Exercise 2.15	
	<p>Write the normal densities in the exponential form (??). What is the natural statistics? What is the natural parametrisation?</p>

2.5.2 Local Regression

2.6 R Complements

2.6.1 Complements: Discretisation

2.6.2 Complements: External Data

2.6.3 Complements: Testing Software

Exercise 2.16	
	For this series of exercises, let $y_i = a + bx_i + \varepsilon_i$ with ε_i iid $\sim N(0, \sigma^2)$ and $x_i = i, i = 1, \dots, 10$.
	Choose a strategy to inspect $lm()$ with regard to the parameter space (a, b, σ^2) . Are there apparent cellular decompositions for the parameters a, b, σ^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.

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

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

2.6.4 R Data Types

2.6.5 Classes and Polymorphic Functions

2.6.6 Extractor Functions

2.7 Statistical Summary

2.8 Literature and Additional References

- [1] Chambers, J.M.; Hastie, T.J. (eds.) (1992): *Statistical Models in S*. Chapman & Hall, New York.
- [5] Jørgensen, B. (1993): *The Theory of Linear Models*. Chapman & Hall, New York.
- [8] R Development Core Team (2004–2008): The R language definition.
- [10] Sawitzki, G. (1994): Numerical Reliability of Data Analysis Systems. *Computational Statistics & Data Analysis* 18.2 (1994), 269–286. <http://www.statlab.uni-heidelberg.de/reports/>.
- [11] 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/>.

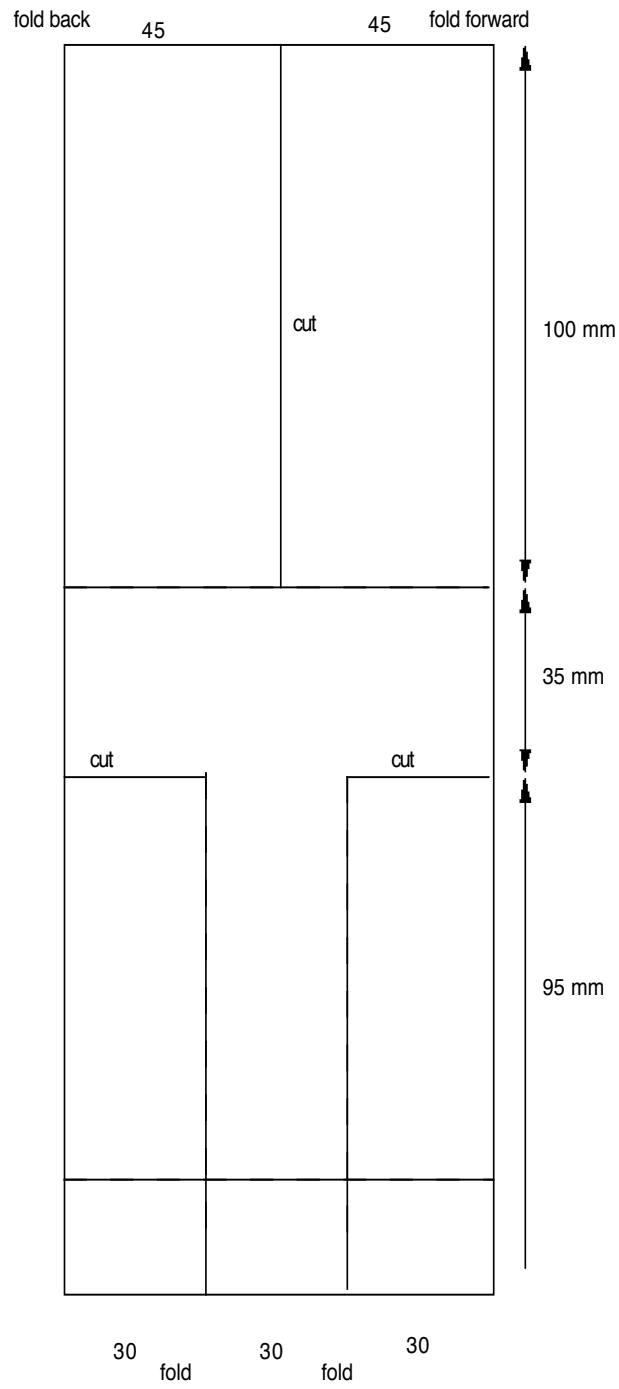


Figure 2.1 *KiwiHopp*

CHAPTER 3

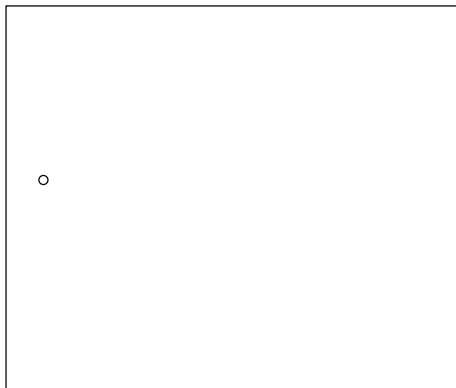
Comparisons

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

Example 3.1: Interactive Location

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

Please click on the circle



Now we wrap up the base function in a timer. We record the coordinates, try to measure the reaction time, and return the results as a list.

Example 3.2: Click Timing

```

----- Input -----
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

```

----- Input -----
dx <- as.data.frame(click1())
dx <- rbind(dx, data.frame(click1()))
dx
----- Output -----
      timestamp      x      y xclick yclick tclick
elapsed 2014-04-23 21:54:22 0.9619 0.6249 0.9609 0.6250 1.944
elapsed1 2014-04-23 21:54:24 0.5477 0.2259 0.7410 0.4125 2.212

```

Exercise 3.1	Click Timing
	<p>Define a function <code>click(runs)</code> that repeats <code>click1()</code> a chosen number <code>runs</code> plus one times and returns the result as a <code>data.frame</code>. The additional first timing should be considered as a “warming up” and is not included in the following evaluations.</p> <p>Select a number <code>runs</code>. Give reasons for your choice of <code>runs</code>. Execute <code>click(runs)</code> and store the result in a file using <code>write.table()</code>.</p> <p>Display the distribution of the component <code>tclick</code> with the methods from Chapter 1 (distribution function, histogram, box-and-whisker plot).</p>

3.1 Shift/Scale Families, and Stochastic Order

Exercise 3.2	Click Comparison
	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> for the right and left hand.
	<p>The recorded data also contain information about the positions. Define a distance measure <i>dist</i> for the deviation. Give reasons for your definition. Perform a right/left comparison for <i>dist</i>.</p> <p>For later analysis, store the results for the right hand and for the left hand in files. named "<i>clickright-xxxx</i>" and "<i>clickleft-xxxx</i>", where <i>xxxx</i> is an identification of you choice. For example, use your initials, the date and some sequential number, such as in "<i>clickright-cs20050416-1</i>".</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 and identically distributed for each treatment. We use the index notation that 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 . The parameter a is called the shift or location parameter.

Speaking in terms of probabilities, the treatment can shift probability mass in quite different ways from 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 (in symbols, $F_1 \prec F_2$), if a variable distributed as F_1 takes rather smaller values than a variable distributed 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 comparison experiment (Exercise 3.2) 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, since the distribution functions are not parallel.

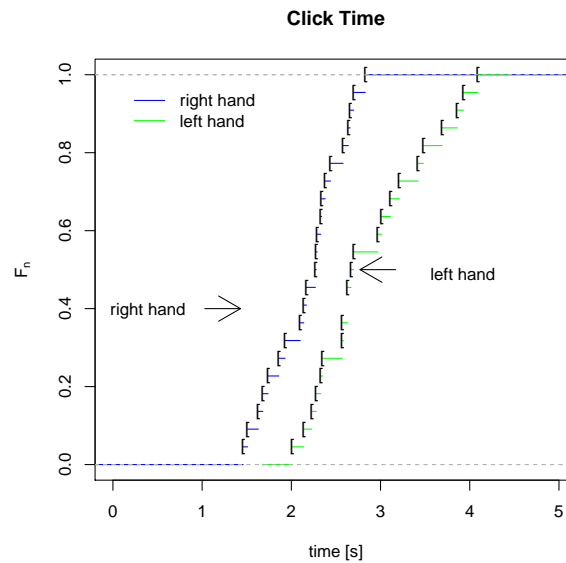


Figure 3.1 *Distribution functions for the right/left click time (samples from one person).*

Exercise 3.3	Stochastic Order
	What does a <i>PP</i> plot for F_1 against F_2 look like if $F_1 \prec F_2$?
	What does a <i>QQ</i> plot for F_1 against F_2 look like if $F_1 \prec F_2$?

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 $N(0, 1)$?

3.2 QQ Plot, PP Plot, and Comparison of Distributions

Exercise 3.5	
	Use the <i>QQ</i> plot to compare the results of the right/left <i>click</i> experiments. Summarise the results.
	Combine the right/left <i>tclick</i> 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 <i>sample()</i> . With <i>par(mfrow = c(2, 2))</i> you arrange the display area so that it shows four plots at a time.
**	For <i>sample()</i> use <i>replace = FALSE</i> . How do you have to apply <i>sample()</i> now to split the joint vector into two vectors with Monte Carlo samples? What differences do you expect in comparison to <i>replace = TRUE</i> ?

Exercise 3.6	
	Find scale and shift parameters for the right/left <i>click</i> data so that, after using these parameters for transformation, the groups match as well as possible. Describe the differences using these parameters. Use a model formulation in terms of a linear model.
	Use the function <i>boxplot()</i> 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.

```
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
}
```

Exercise 3.7	Scale/Shift Standardisation
	<p>This algorithm is only appropriate for symmetric distributions.</p> <p>Combine it with a power transformation as in Section 2.5 (page 23) to symmetrise a distribution and give an algorithm that can be applied to asymmetric transformations.</p>

Exercise 3.8	Two-Sample Monte Carlo Bands
*	<p>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.)</p> <p>For the bands, you can use an overlay of line plots.</p> <p><i>Hint:</i> Use the function <code>sample()</code> to generate random permutations.</p>

Exercise 3.9	
**	<p>Augment the <i>PP</i> plot and <i>QQ</i> plot for the <code>click</code> experiments by permutation bands that cover 95% of the permutations.</p>
**	<p>Generate new plots from the <i>PP</i> plots and <i>QQ</i> plots by adding Monte Carlo bands from permutations. Use an envelope of 19 Monte Carlo samples.</p> <p><i>Hint:</i> Use function <code>sample()</code> to draw a random sample of sample size n_1 from $x = (Y_{11}, \dots, Y_{1n_1}, Y_{12}, \dots, Y_{1n_2})$.</p>
	<p><i>Hint:</i> See <code>help(sample)</code>.</p>

Exercise 3.10	
*	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.

Exercise 3.11	
*	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

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 the 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 π 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 $(Y'_j)_{j=1, \dots, m}$.

Exercise 3.12	
*	Implement this algorithm and enhance the <i>PP</i> plot by adding simulated <i>PP</i> plots generated by 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 rejection, i.e., an error. Inspect the distribution of the error level under the simulated conditions.

Exercise 3.13	
	Use the <i>QQ</i> plot for a pair-wise comparison of the results of the helicopter experiment from Chapter 2. Summarise your results.

Exercise 3.14	
	Inspect the implementation of <code>qqnorm()</code> . Implement an analogous function for the <i>PP</i> plot and apply it to the helicopter data.

3.3 Tests for Shift Alternatives

Exercise 3.15	
*	Use a simulation to inspect the distribution of \bar{Y} , $\widehat{Var}(Y)$ and the t statistic for Y from a uniform distribution $U[0, 1]$ with sample size $n = 1, \dots, 10$. Compare the distributions from the simulation with the corresponding normal, χ^2 resp. t distribution.
	Use a simulation to inspect the distribution of \bar{Y} , $\widehat{Var}(Y)$ and the t statistic for Y from a mixture, consisting at 90% from an $N(0, 1)$ - and at 10% from an $N(0, 10)$ distribution, with sample size $n = 1, \dots, 10$. Compare the distributions from the simulation with the corresponding normal, χ^2 resp. t distribution.

Exercise 3.16	
	Use the Wilcoxon test to compare the results of the right/left <i>click</i> experiment. Use both variants, the approximative test <code>wilcox.test()</code> and the exact Wilcoxon test <code>wilcox_test()</code> .

Exercise 3.17	Click Project
***	<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"> • The response time comprises reaction time, time for the large scale movement of the mouse, time for fine adjustment, etc. • 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. • Subsequent records may be affected by a training effect, or by a tiring effect. <p>Can you modify the experiment or the evaluation so that differences in the reaction time components can be investigated? 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>

Exercise 3.18	Power Comparison
	<p>Use the shift/scale families of $N(0, 1)$ and $t(3)$ and design a setting to compare the performance of the Wilcoxon test with that of the t-test for each of these families.</p> <p>Perform the comparison in a simulation with sample sizes $n_1 = n_2 = 10, 20, 50, 100$ and summarise your results.</p>
	<p>Do an analogous comparison using simulation data from the log-normal distribution.</p>

3.4 A Road Map

3.5 Power and Confidence

3.5.1 Theoretical Power and Confidence

3.5.2 Simulated Power and Confidence

3.5.3 Quantile Estimation

Exercise 3.19	Interquartile Interval
	Write a function that calculates the coverage probability $n \mapsto P(\text{med}(X) \in I_Q)$ where I_Q is the empirical interquartile interval. <i>Hint:</i> Use (??).
	What is the minimal sample size so that the interquartile box covers the median with at least 90% confidence?

3.6 Qualitative Features of Distributions

3.7 Statistical Summary

3.8 Literature and Additional References

[14] William N. Venables and Brian D. Ripley, B (2002): *Modern Applied Statistics with S*.

Springer, Heidelberg.

[13] William N. Venables, W.N.; and Brian D. Ripley (2000): *S Programming*.

Springer, Heidelberg.

[7] Rupert G. Miller (1981): *Simultaneous Statistical Inference*.

Springer, Heidelberg.

CHAPTER 4

Dimensions 1, 2, 3, ..., ∞

4.1 R Complements

4.2 Dimensions

4.3 Selections

4.4 Projections

4.4.1 Marginal Distributions and Scatter Plot Matrices

Exercise 4.1	
	Generate a scatterplot matrix for the diabetes data set that shows a histogram of the variables in the diagonal panels. <i>Hint: See <code>help(pairs)</code>.</i>

4.4.2 Projection Pursuit

Input

```

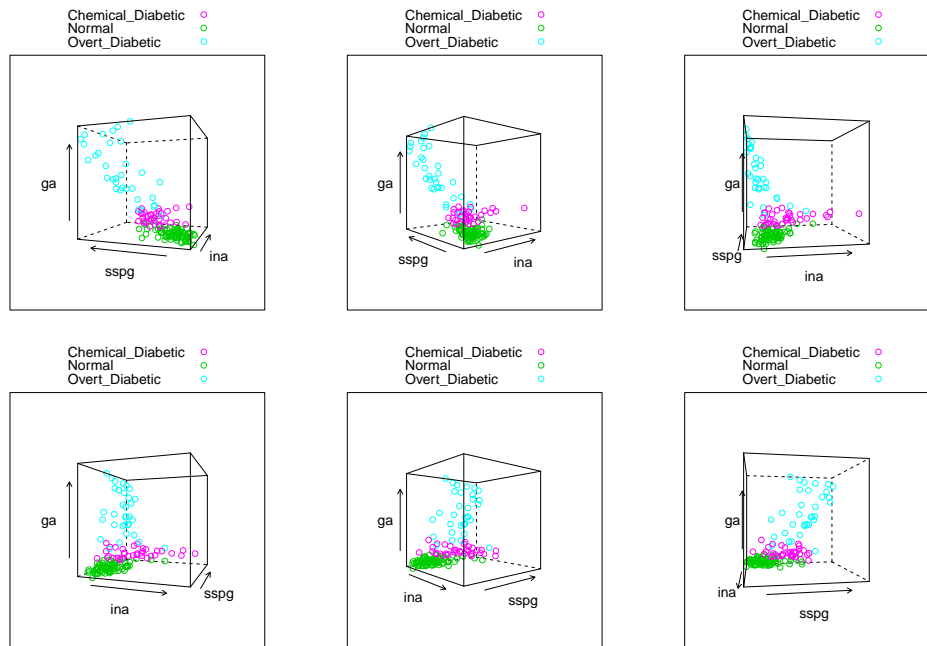
library("lattice")
diabcloud <- function(y, where, more = TRUE, ...) {
  print(cloud(ga ~ ina + sspg, data = chemdiab, groups = cc,
    screen = list(x = -90, y = y), distance = .4, zoom = .6,
    auto.key = TRUE, ...),
    split = c(where, 3, 2), more = more)
}
supsym <- trellis.par.get("superpose.symbol")
supsymold <- supsym
supsym$col = c("magenta", "green3", "cyan")
trellis.par.set("superpose.symbol" = supsym)
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, supsymold, supsym)

```



See Colour Figure ??.

Exercise 4.2	
	<p>Modify this example so that you get an impression of the three-dimensional structure. Try to use an animated sequence. You can use <code>sys.wait()</code> if it is available on your system to control the time sequence, or use <code>devAskNewPage()</code> to give interactive control for new images.</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>

4.4.3 Projections for Dimensions 1, 2, 3, ... 7

4.4.4 Parallel Coordinates

Exercise 4.3	
	For the <code>chemdiab</code> data set, prepare a (written!) report about the relation between the variables that you can recognise in the parallel coordinate plot.
	<p>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.</p> <p>Which sequence of variables gives the simpler display?</p> <p>Which relations between the variables are visible in both?</p> <p>Which relations appear only in one of the arrangements?</p>

4.5 Sections, Conditional Distributions and Coplots

From an abstract point of view, sections are conditional distributions of the type $P(\cdot | X = x)$. But they are only reliable where the section defines a condition that has positive probability. To make the idea of restricting the view on conditional distributions applicable to data, we thicken the sections. Instead of considering conditional distributions of the type $P(\cdot | X = x)$ we consider $P(\cdot | \|X - x\| < \varepsilon)$, where ε 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 [2]. 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.

As a first tool, R provides the possibility to 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 panel the scatterplot of two variables, given the condition.

The coplot can be inspected for patterns. 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 at 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 may 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,
       axlabels = function(f) abbreviate(levels(f)),
       number = 6, overlap = 0.5, xlim, ylim, ...)
co.intervals(x, number = 6, overlap = 0.5)
```

Arguments

formula	a formula describing the form of conditioning plot. A formula of the form $y \sim x \mid a$ indicates that plots of y versus x should be produced conditional on the variable a . A formula of the form $y \sim x \mid a * b$ 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.
data	a data frame containing values for any variables in the formula. By default the environment where <code>coplot</code> was called from is used.
given.values	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.

<code>panel</code>	a function(<code>x</code> , <code>y</code> , <code>col</code> , <code>pch</code> , ...) which gives the action to be carried out in each panel of the display. The default is <code>points</code> .
<code>rows</code>	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.
<code>columns</code>	the number of columns in the panel layout array.
<code>show.given</code>	logical (possibly of length 2 for 2 conditioning variables): should conditioning plots be shown for the corresponding conditioning variables (default <code>TRUE</code>).
<code>col</code>	a vector of colors to be used to plot the points. If too short, the values are recycled.
<code>pch</code>	a vector of plotting symbols or characters. If too short, the values are recycled.
<code>bar.bg</code>	a named vector with components <code>"num"</code> and <code>"fac"</code> giving the background colors for the (shingle) bars, for numeric and factor conditioning variables respectively .
<code>xlab</code>	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.
<code>ylab</code>	character; labels to use for the y axis and any second conditioning variable.
<code>subscripts</code>	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.
<code>axlabels</code>	function for creating axis (tick) labels when x or y are factors.
<code>number</code>	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 x and y direction. When <code>overlap < 0</code>, there will be <i>gaps</i> between the data slices.
<code>xlim</code>	the range for the x axis.
<code>ylim</code>	the range for the y 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` \geq `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 $(\text{number} \times 2)$ 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 sensible, but 'show off':
  coplot(Life.Exp ~ state.region | Income * state.division,
```

```
    panel = panel.smooth)  
  })
```

Exercise 4.4	Earthquakes
	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 earthquake? (You may have to choose a different model formula for the plots.)

4.6 Transformations and Dimension Reduction

Exercise 4.5	Iris Classification
	Use the methods from Section 4.4 and 4.5 to inspect the data set. Can you see classification rules that give a classification of the three species to a large extent?

4.7 Higher Dimensions

4.7.1 Linear Case

Partial Residuals and Added Variable Plots

Exercise 4.6	Added Variables
	<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 produces a scatterplot matrix of these residuals.</p> <pre data-bbox="478 1032 970 1149">pairslm <- function(model, x, ...) { x\$fit <- lm(model, x)\$fit pairs(x, ...) }</pre> <p>Add title, legends, etc.</p> <p>Use the “trees” data set as an example.</p>

4.7.2 Non-Linear Case

Example: Cusp Non-Linearity

4.7.3 Case Study: Melbourne Temperature Data

4.7.4 Curse of Dimensionality

4.7.5 Case Study: Body Fat

Exercise 4.7	
	<p>Use functions <code>plot()</code>, <code>identify()</code>, and <code>text.id()</code> to generate the following output:</p>

Exercise 4.8	
*	<p>Use</p> <pre>library(leaps) lm.reg <- regsubsets(body.fat ~ age + BMI + neck + chest + abdomen + hip + thigh + knee + ankle + bicep + forearm + wrist + weightkg + heightcm, data = fat)</pre> <p>and inspect the result with</p> <pre>summary(lm.reg) plot(lm.reg, scale = "r2") plot(lm.reg, scale = "bic") plot(lm.reg, scale = "Cp")</pre> <p><i>Hint: See <code>help(plot.regsubsets)</code>.</i></p>
*	Use the function <code>leaps()</code> for model selection.

Exercise 4.9	
	Remove the obvious outliers and rearrange the variables starting from the body volume so that on average the correlation between subsequent variables is maximized.

Exercise 4.10 Think !	
	Draw a sketch of a member doll that shows which body geometry features are represented by the next principal component <i>PC4</i> , ..., <i>PC10</i> . For a start, you can concentrate on the signs of the variable weights.

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 <i>age</i> in the model. How exactly do you include <i>age</i> in the model?
**	The function <i>mvr()</i> in <i>library(pls)</i> [15] is available to perform a regression based on principal components. Use this function for regression. What is the difference between this estimation and the 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:

```
fat.eval <- fat[fat$train == FALSE, ] Input
pred <- predict.lm(lm.volf, fat.eval, se.fit = TRUE)
```

Exercise 4.12	
*	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.

48

4.8 High Dimensions

4.9 Statistical Summary

References

- [1] John M. Chambers and Trevor J. Hastie. *Statistical Models in S*. Chapman & Hall, London, 1992.
- [2] George W. Furnas and Andreas Buja. Prosection views: Dimensional inference through sections and projections. *J. Comput. Graph. Statist.*, 3(4):323–385, 1994.
- [3] Peter Gänßler and Winfried Stute. *Wahrscheinlichkeitstheorie*. Springer, Heidelberg, 1977.
- [4] Robert Gentleman and Ross Ihaka. Lexical scope and statistical computing. *Journal of Computational and Graphical Statistics*, 9:491–508, 2000.
- [5] Bent Jørgensen. *The Theory of Linear Models*. Chapman & Hall, New York, 1993.
- [6] Paul Massart. The tight constant in the Dvoretzky-Kiefer-Wolfowitz inequality. *The Annals of Probability*, 18(3):1269–1283, July 1990.
- [7] Rupert G. Miller. *Simultaneous Statistical Inference*. Springer, New York, 1981.
- [8] R Development Core Team. *The R language definition*, 2008.
- [9] R Development Core Team. *Writing R Extensions*, 2008.
- [10] Günther Sawitzki. Numerical reliability of data analysis systems. *Computational Statistics & Data Analysis*, 18(2):269–286, 1994.
- [11] Günther Sawitzki. Report on the numerical reliability of data analysis systems. *Computational Statistics & Data Analysis*, 18(2):289–301, 1994.
- [12] Galen R. Shorack and Jon A. Wellner. *Empirical Processes with Applications to Statistics*. Wiley, New York, 1986.
- [13] William N. Venables and Brian D. Ripley. *S Programming*. Statistics and Computing. Springer, New York, 2000.
- [14] William N. Venables and Brian D. Ripley. *Modern Applied Statistics with S*. Springer, Heidelberg, fourth edition, 2002.
- [15] Ron Wehrens and Bjørn-Helge Mevik. *pls: Partial Least Squares Regression (PLSR) and Principal Component Regression (PCR)*, 2007. R package version 2.1-0.

Subject Index

- conditioned, **39**
- data features, 36
- data types, 25
- distribution
 - Gaussian, 31
- distribution function, 2
- factor, 17
- function, **12**
- histogram, 2, 6
- interactive, 27, 38
- lattice, 37
- model
 - formula, *see* Wilkinson-Rogers notation
 - generalised linear, 23
 - linear, 17, 19, 29
- moment, 8
- Monte Carlo, 3, 20, 31
- plot
 - box-and-whisker, 11
 - coplot, 39
 - histogram, 2, 6
 - PP*, 31
 - QQ*, 31
 - Tukey-Anscombe, 20
- polymorphic, 25
- power, 35
- quantile, 8, 35
- quartile, 31
- regression
 - non-linear, 23
 - principal component, 47
 - robust, 32
 - samples
 - Monte Carlo, 31
 - shift family, **29**
 - simulation, 35
 - stochastically smaller, **29**
- Wilkinson-Rogers notation, 19