

Computational Statistics

An Introduction to 

Günther Sawitzki



CRC Press
Taylor & Francis Group

A CHAPMAN & HALL BOOK

Computational Statistics

An Introduction to

Exercises

Günther Sawitzki
StatLab Heidelberg
April 28, 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
References	17
Functions and Variables by Topic	19
Function and Variable Index	19
Subject Index	19

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><code>par(mfrow = c(2, 3))</code></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>plot(sin(1:100))</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>hist(runif(100), breaks = c(0, 0.5, 1))</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

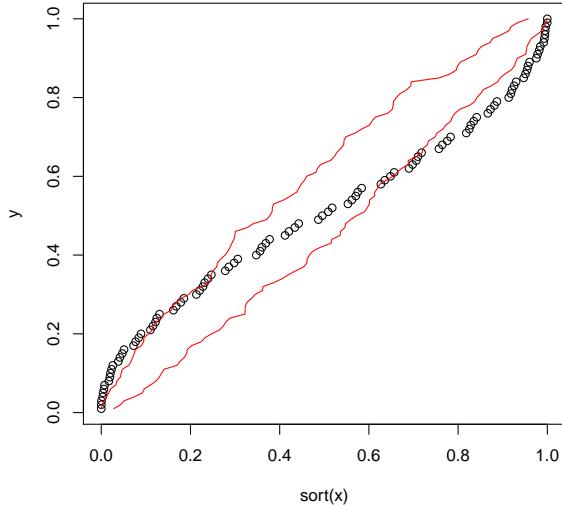
Exercise 1.5	
	Using <code>help(rbeta)</code> you get information about the functions available to work with beta distributions. Generate plots for the densities of the beta distribution for $n = 16, 32, 64, 128$ and $i = n/4, n/2, 3n/4$. Use the function <code>curve()</code> to generate the plots. For 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 `help()`-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.)	
apply()	applies a function to the rows or columns of a matrix. <i>Example:</i> <code>samples <- apply(samples, 2, sort)</code> sorts by column.
outer()	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 <code>nrsamples = 19</code>.</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 <code>sim</code> 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 <code>sim</code> 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. [3], 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	
	<p>Using <code>help(ks.test)</code> you get information on 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 runif(100) sin(1:100) rnorm(100)</pre> <p>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.</p>

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

Exercise 1.12	
	<p>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:</p> <pre>(3 3 3 3 3) (1 2 5 3 3) (0 0 9 0 6).</pre>

Exercise 1.13	
	<p>Which results do you expect if you use a χ^2 test to check the following vectors for a uniform distribution?</p> <pre>1:100 runif(100)</pre> <p>(cont.)→</p>

Exercise 1.13	(cont.)
	<pre>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 component of its result.</p>

Exercise 1.14	
*	<p>Sketch comparable test environments for fixed and adaptive choice of histogram cells.</p> <p>For fixed and for adaptive choice of histogram cells draw $s = 1000$ samples of size 50 from <code>rnorm()</code>. Calculate in both settings the formal χ^2 statistics and plot its distribution functions.</p> <p>Compare the distribution functions.</p>

Exercise 1.15	
	<p>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>You have to choose a χ^2 test. 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 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>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:</i> <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.</p> <p>Using</p> $\text{contour}()$ <p>you can generate a contour plot.</p> <p>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) $x \quad x = 1..100 \bmod m$ for convenient m ii) $\sin(x) \quad x = 1..100$ 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>

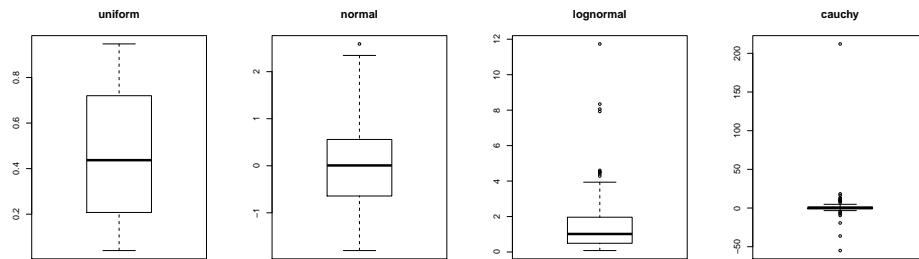
1.4 Moments and Quantiles

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

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-and-whisker plot to show this interval.

(cont.)→

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 <i>apply()</i> to avoid the loop. <i>Hint:</i> See Example 1.1 (page 3).

Exercise 1.26	
	Use <i>rnorm()</i> 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 <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

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>eecdfe</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									
	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>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" style="margin-left: auto; margin-right: auto;"> <thead> <tr> <th>Generator</th><th>Features</th></tr> </thead> <tbody> <tr> <td>uniform</td><td>flat density</td></tr> <tr> <td>lognormal</td><td>unimodal, skewed, long tail</td></tr> <tr> <td>Cauchy</td><td>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	
	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.
	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.

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 <code>x</code>, the statement <code>prompt(x)</code> generates a skeleton upon which you can build a documentation for <code>x</code>.</p>

Compilation

1.7 Statistical Summary

1.8 Literature and Additional References

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

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

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

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

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

References

- [1] Peter Gänßler and Winfried Stute. *Wahrscheinlichkeitstheorie*. Springer, Heidelberg, 1977.
- [2] Robert Gentleman and Ross Ihaka. Lexical scope and statistical computing. *Journal of Computational and Graphical Statistics*, 9:491–508, 2000.
- [3] Paul Massart. The tight constant in the Dvoretzky-Kiefer-Wolfowitz inequality. *The Annals of Probability*, 18(3):1269–1283, July 1990.
- [4] R Development Core Team. *Writing R Extensions*, 2008.
- [5] Galen R. Shorack and Jon A. Wellner. *Empirical Processes with Applications to Statistics*. Wiley, New York, 1986.

Subject Index

distribution function, 2

function, **12**

histogram, 2, 6

moment, 8

Monte Carlo, 3

plot

 box-and-whisker, 11

 histogram, 2, 6

quantile, 8