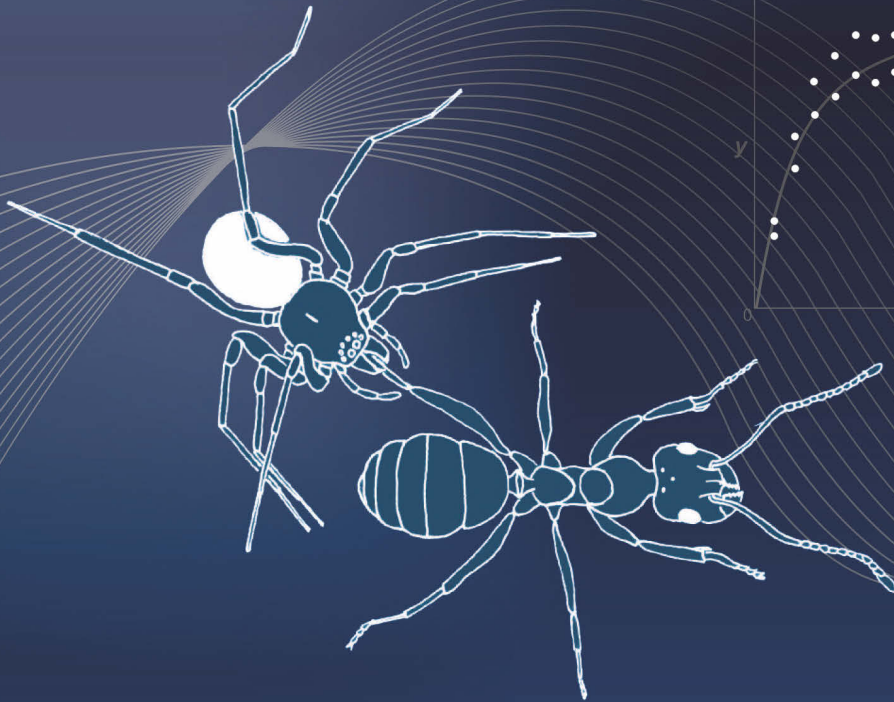


MODERN ANALYSIS OF BIOLOGICAL DATA

GENERALIZED LINEAR
MODELS IN **R**



STANO PEKÁR
MAREK BRABEC

muni
PRESS

MODERN ANALYSIS OF BIOLOGICAL DATA
GENERALIZED LINEAR MODELS IN R

STANO PEKÁR, MAREK BRABEC

MODERN
ANALYSIS
OF BIOLOGICAL DATA
GENERALIZED LINEAR MODELS
IN R

STANO PEKÁR
MAREK BRABEC

Masaryk University, Brno 2016

http://www.muni.cz/press/books/pekar_en

Pekár S. & Brabec M. 2016. *Modern Analysis of Biological Data: Generalized Linear Models in R*. Masaryk University Press, Brno.

This book was supported by Masaryk University Project No. MUNI/FR/1304/2014.

Text © 2016 Stano Pekár, Marek Brabec
Illustrations © 2016 Stano Pekár
Design © 2016 Ivo Pecl, Stano Pekár, Grafique
© 2016 Masarykova univerzita

ISBN 978-80-210-8106-2 (online : pdf)
ISBN 978-80-210-8019-5 (Paperback)

Foreword

1 Introduction

1.1	How to read the book	3
1.2	Types of variables	5
1.3	Conventions	6

2 Statistical software

2.1	The R Environment	7
2.2	Installation and use of R	9
2.3	Basic operations	11
2.4	Data frames	18

3 Exploratory data analysis (EDA)

3.1	Expected value	23
3.2	Variance	25
3.3	Confidence intervals	26
3.4	Summary tables	27
3.5	Plots	28
3.5.1	Distribution plots	32
3.5.2	Scatter plots	35
3.5.3	Box plots	35
3.5.4	Lattice plots	37
3.5.5	Interaction plots	38
3.5.6	Bar plots	39
3.5.7	Paired plots	40
3.5.8	3D plots	40
3.5.9	Plots with whiskers	40
3.5.10	Curves	41

4 Statistical modelling	
4.1 Regression model	43
4.2 General linear model	45
4.3 Generalized linear model	47
4.4 Searching for the “correct” model	51
4.5 Model selection	53
4.6 Model diagnosis	54
5 The first trial	
5.1 An example	61
5.2 EDA	61
5.3 Presumed model	63
5.4 Statistical analysis	63
5.4.1 ANOVA table of Type I	65
5.4.2 Nonlinear trends	67
5.4.3 Removal of model terms	70
5.4.4 Comparison of levels using contrasts	74
5.4.5 Contrasts and the model parameterization	77
5.4.6 Posterior simplification	83
5.4.7 Diagnosis of the final model	85
5.5 Conclusion	88
6 Systematic part	
6.1 Regression	90
6.2 ANOVA and ANODEV	93
6.3 ANCOVA and ANCODEV	94
6.4 Syntax of the systematic part	96
7 Random part	
7.1 Continuous measurements	100
7.2 Counts and frequencies	102
7.3 Relative frequencies	104
8 Gaussian distribution	
8.1 Description of LM and GLM	107
8.2 Regression	108
8.3 Weighted regression	116
8.4 Multiple regression	120

8.5	Two-way ANOVA	132
8.6	One-way ANCOVA.....	141
9 Gamma and lognormal distributions		
9.1	Description of the Gamma model.....	147
9.2	Description of the lognormal model.....	148
9.3	Regression.....	149
9.4	Two-way ANODEV.....	156
9.5	Two-way ANCOVA.....	163
10 Poisson distribution		
10.1	Description of the Poisson model.....	169
10.2	One-way ANODEV.....	170
10.3	Overdispersion and underdispersion	175
10.4	Multiple regression	176
10.5	One-way ANCODEV	183
10.6	Three-way ANODEV (Contingency table)	190
11 Negative-binomial distribution		
11.1	Description of the negative-binomial model.....	199
11.2	One-way ANODEV.....	200
12 Binomial distribution		
12.1	Description of binomial model.....	210
12.2	Two-way ANODEV.....	212
12.3	Overdispersion and underdispersion	218
12.4	Regression.....	219
12.5	One-way ANCODEV	226
12.6	Binary one-way ANCODEV	231
References		
Index		
	Subject index.....	239
	R functions and their arguments.....	243

FOREWORD

This book is meant especially for students and scholars of biology, i.e. biologists who work in natural science, at agricultural, veterinary, pharmaceutical and medical faculties or at research institutes of a similar orientation. It has been written for people who have only a basic knowledge of statistics (for example, people who have attended only a Basic statistics/Biostatistics course) but who need to correctly analyse the data resulting from their observations or experiments.

The generally negative attitude of biologists towards mathematics is well known. It is precisely why we have tried to write the book in a relatively simple style – with minimal mathematical requirements. Sometimes, the task turned out to be easy, other times not that easy, and sometimes it became almost impossible. That is why there are still some mathematical equations in almost all chapters of the book (even though they are used in a simplified form in order to be more apprehensible to less experienced readers). Despite this fact, the book includes much less mathematical and statistical theories than it is common for standard statistical literature.

The book is mainly built on examples of real data analyses. They are presented from the very beginning to the end, from a description and determination of objectives and assumptions to study conclusions. They thus simulate (even though in a simplified way) the procedure usually used when preparing a paper for a scientific journal. We believe that practical experience with data analyses is irreplaceable. Because of the anticipated biology-oriented readers, we selected examples from the areas of ecology, ethology, toxicology, physiology, zoology and agricultural production. All these data were analysed during various previous research projects. They have been adjusted in order to suit the pedagogical intentions of this book. For example, the original long and complex Latin names of species have been replaced with a generic short name (e.g., specA).

Finally, we would like to thank all our colleagues without whose help this book would never have been written. First of all, we would like to thank to Vojtěch Jarošík (in memoriam), for introducing GLM to the first author of the book during his studies at the university, thus igniting his interest in statistics generally; Alois Honěk for many consultations, and our colleagues from the Crop Research Institute in Prague-Ruzyně and the students of the Faculty of Science of the Masaryk University in Brno for inspiring comments to the original text of the book. Finally, we would also like to thank the following colleagues of ours who have kindly let us use their (though adjusted) data for presenting examples in this book: T. Bilde, A. Honěk, J. Hubert, D. Chmelař, J. Lipavský, M. Řezáč, P. Saska and V. Stejskal.

We welcome any comments regarding the text and content of the book. Please direct them to the following email addresses: pekar@sci.muni.cz and/or mbrabec@cs.cas.cz.

December 2015

Stano Pekár
Marek Brabec

Let us start with a demonstration of two standard situations. The first one took place after a thesis defence when a student complained to another student: “Supposedly I used the wrong statistical test.” Other situations can occur, for example, in a hallway of a research institute where a biologist reproaches a colleague for the way he/she presented his/her results: “The data should be analysed more properly.” Both situations have one thing in common – a desperate reference to the statistics. Indeed, statistical data analysis forms an integral part of scientific publications in many biological (and other) fields, thus accompanying bio-logists throughout their careers. In some fields, such as taxonomy, statistical analyses may play just a marginal role. For other fields, such as ecology or physiology, it is often almost a cornerstone of many new findings. In these fields, shortcomings in statistical analysis can have catastrophic consequences. It can easily happen that a report on a good experimental study (e.g., in the form of a scientific paper) will not be complete without a statistical analysis, in which case the results of the study can become completely useless and all the previous effort of its authors can thus be wasted.

The only way to prevent such disasters is to strive to understand the statistics or, at least, to find somebody who understands it. Obviously, conducting practical data analyses are much easier today than ever before. This is due to the development of personal computers and subsequent developments in computational algorithms, which have literally meant a revolution for data analyses. Fifty years ago, even a simple statistical analysis, using a calculator and a pen, would take several hours and sometimes even days (while, at the same time, it was not easy to avoid calculation errors etc.). Today, using computers, even a relatively complicated analysis can take less than a minute or even just a few milliseconds. Preparing a plot is often easier and faster than, for example, preparing a coffee. However, technical improvements have led to increased demands for using adequate statistical methods. While simpler statistical methods were preferred in the past, despite the fact they were not the most suitable, today the emphasis is put on using methods that correspond in the relevant aspects to the actual data at hand as closely as possible. This is because the computing requirements do not represent an insurmountable obstacle any more. Unlike the former “universal” simple procedures, application of statistics today utilises such methods and models that realistically consider the important characteristics of the data and of the studied processes. Very often this means that a given method or a statistical model can and should be adjusted to the real situation at hand. In short, models should be adjusted to the data and not vice versa!

Nevertheless, it is obviously not always easy to comply with this requirement. In fact, creative and useful practical analyses need theoretical knowledge about numerous models and methods as a pre-requisite. Moreover, certain experience with practical data analyses, with

the application of various models on real data, with model building as well as estimation procedures, with conducting appropriate tests, etc. is also necessary. It is clear that one cannot become an experienced data analyst only from reading books. To get such experience, you have to put in some work and, most of all, some thinking. Guidance and examples can assist you in this process.

This book attempts to help exactly along these lines by presenting examples of particular analyses, including the specification of the problem of interest, development of a statistical model and formulation of possible conclusions. The book is not, and does not even aspire to be, a manual for selecting the “best” method for a particular data analysis (it is our strong opinion that, for various reasons, such a manual cannot be ever made). Instead, the book tries to demonstrate how to think about particular statistical models, how to use them and also how not to use them – to point out many of the problems and errors that can occur (and do indeed occur in practical analyses and even in published papers). We demonstrate various general approaches on particular (hopefully biologist-engaging) examples and we also present their detailed implementation in the R language, thus allowing everybody to try them for themselves using the data provided in the book as well as their own datasets of a similar nature.

Since regression is a very powerful instrument, used very often in biological studies, this book is almost exclusively dedicated to regression models. It assists in solving important questions of the following type: how does variable y depend on variable x or how can we predict the value of a new observation when we know the value of one or several so-called explanatory variables. However, we will talk about regression in a somewhat wider context compared to what you know from basic statistics courses. Namely, we will use the concept of the so-called Generalized Linear Models (GLM). Their name explicitly emphasises that they are generalized linear regression models – that is, generalisations of models that many people know under the term of general linear model. The GLM generalization is very useful from a practical point of view since it allows correct analyses of various data that cannot be processed using a standard (linear) regression without (very) rough approximations and (crude) simplifications. GLM represents a relatively large class of models with a unified statistical theory and very universal computational implementations. It is a class by the means of which you can analyse a wide range of data, such as weight measurements, concentration, number of individuals, or relative frequency of a phenomenon (i.e. various continuous unrestricted, continuous positive, discrete positive data).

From a formal point of view, we will use only univariate GLM models (or statistical methods based on them), and only those that are suitable exclusively for independent measurements (observations). In order to analyse data with some kind of correlation (for example, in the case of repeated measurements on the same individuals), somewhat different and more complex models and methods need to be used. We address these in another book (Pekár & Brabec 2012).

As we have already stated, the objective of this book is to assist practical users of statistics in the process of formulating and using statistical models – and thus also in selecting suit-

able methods for analyses of various data types. At the same time, we would like to motivate them to seek assistance from a professional statistician if they find that the complexity of the model, study design and/or data complexity exceed their abilities and experience with statistical modelling. Things are not always as easy as they may appear at first. The same data can be analysed using various methods, depending, for example, on the objectives of a particular analyst. In fact, data include a large amount of information but we are often interested only in extracting just some part of it at a given moment. That is why the method selected depends on what we regard as salient features of the data (and what we are willing to leave out), as well as what (and for what purposes) we want to extract from the observations/measurements. Nevertheless, it is often the case that multiple procedures (with different characteristics) can be used for analysing even a single aspect. In this book, we will mostly present only one procedure without denying the existence of other approaches. Otherwise, the book would become significantly longer.

We have applied a similar approach when describing the R language and discussing the use of particular objects. The R language is, in fact, very rich. Like in any other programming language environment, the same result can be often obtained in several ways. There is not enough room here to list them all (if you are interested, you can consult the corresponding manuals, e.g. Zoonekynd 2007). We chose those that we consider to be the easiest and most practical for a beginner.

1.1 How to read the book

The text of this book combines examples of selected statistical methods and descriptions of the R language as a statistical environment. Both are then illustrated on practical data analyses.

How to read this book? It depends on your knowledge and experience with statistical analyses as well as with the R environment. Those of you who have not done any (or almost any) data analysis since they attended their basic statistics/biostatistics course and who have never worked with R before should read the book from the beginning to the end. You who already know the environment R (at least the basics) and have been using regression in your work, can read the book somewhat non-systematically – picking the chapter that deals with data and/or methods you are currently interested in.

The most important part of the first chapter is the section 1.2, which defines variables. You certainly should not skip this part even if you believe that you are clear about distinguishing among variable types. This is because our definitions may differ from the definitions you may have encountered in other books and/or courses.

Chapter 2 describes the installation and use of the R software, which we will be using in this book for various data analyses. If you have never heard about R (or you have heard about it but never actually worked with it), this chapter is here especially for you. Apart from the installation of the software, you will also learn both some general principles for working in the R environment and important commands used repeatedly for various data manipulations later

in this book. Moreover, the chapter explains how to enter data into the R environment. All of the above is absolutely essential for being able to use the software successfully, not only for practicing the examples of this book, but also for practical data analyses of your own.

Chapter 3 shows some of the methods of so-called exploratory data analysis. This analysis means calculation of the basic statistical characteristics of the examined data sets as well as their informal comparison using plots, tables and other instruments. You will find there a general description of several useful R commands, by means of which you can create the majority of plots, with which we will work in later chapters.

The next four chapters (Chapters 4-8) are oriented rather more generally and even include some theory. Do not be misled and do not skip them thinking that they are not important from a practical point of view! In fact, they are some of the most important ones. This is because you can use their content for data analyses even if you use different software or, for example, for general considerations related to the strategy of statistical analyses. Particularly, Chapter 4 addresses the issue of how to work with statistical models of the type that will be used later in the book. Model formulation (traditionally in the form of a mathematical formula) forms a steppingstone, on which statistical analysis is based. In this chapter, we will talk about regression models and we will discuss what GLM actually is.

Chapter 5 presents the first example of a concrete and non-trivial data analysis. The presented example is not as simple as motivational examples usually are. It is more demanding, which allows us to discuss more of the aspects you will encounter when working with other examples (and in practical analyses of your own). The analysis is described in detail and commented abundantly (with references to the basic rules, which analysts should observe). In addition to the previous, some R outputs are described and interpreted and important definitions are stated (for example, the definition of contrasts).

Chapter 6 then goes deeper and deals with various systematic parts of a GLM model. It describes basic model classification based on the character of the explanatory variables. However, most importantly, it demonstrates in general terms how to translate a mathematical model into the R language and suggests how to interpret the model after the analysis is done.

Chapter 7 can be used as a simple “key”, based on which readers (beginners) can try to decide which particular method to use (you will see yourself that, once you acquire more knowledge and experience, this decision will require even more thinking). We assume that the reader will be coming back to the book to use it as an aid for his/her own data analyses. If the reader remembers the general knowledge from Chapters 1-6, he/she can go straight to Chapter 7, which will direct him/her to a chapter where he/she will find an analysis of an example similar to his/her own; however, if you cannot decide which method to use, continue reading subsequent chapters in the order they are presented. Either way, we certainly recommend that you initially read Chapters 8-12 completely.

Chapters 8-12 are based on several examples, but they explain also theory, which you will encounter during these (and other) analyses. In general, the analyses of all examples from

these chapters follow a similar plan. The chapters are written in a way that makes them mutually independent (they can be read separately). As a result of this, you encounter similar problems repeatedly. When you pick a chapter, you should always read it as a whole. Do not try to follow it merely as a concrete report of a particular analysis.

1.2 Types of variables

Definitions and a detailed description of the characteristics of various variable types form the topic of a basic statistical course, which we certainly do not intend to repeat here. Let us just remind you of a few important facts that will be the most useful later in the book.

There are several possible viewpoints based on which variables can be classified. For our purposes, we will only need the following variable types:

Response variable: a (random) variable, the variability of which we attempt to explain by means of a statistical (regression) model. In other words, it is a variable that we want to model using a single explanatory variable or multiple explanatory variables. In this book we will exclusively address univariate models. That is, we will always model only one random response variable at a time (in contrast to multivariate models with several response variables considered in one model simultaneously). Depending on the particular GLM model type, our response variable will be either continuous or discrete, but always numeric (i.e. the values are numbers).

Explanatory variable: a variable by the means of which we explain the values of a response variable. Even in the case of univariate models, a response variable can be modelled by a single explanatory variable or by multiple explanatory variables. These can be either numeric or categorical (the values are characters and character chains that correspond to a given code marking groups/categories). Numerical variables can be continuous or discrete. We will mark continuous explanatory variables using lowercase letters, for example x . We will mark their value for the i th observation using index i (e.g. x_i). We will call continuous variables **covariates**. We will mark **categorical** variables using uppercase letters (for example A). a categorical variable can always take at least two different values, which we call **levels** (for example, “male“ and “female“). Then the j th level of such a variable is marked with index j , for example as A_j , a categorical explanatory variable (with categories denoted by characters or numbers) will be also called a **factor**, as in the analysis of variance (ANOVA).

Weights represent a special case of variable. They determine relative weights of individual observations within a given data set. By default, functions in R use the same (unitary) weights for all observations. Externally entered weights are useful when a scheme with equal weights is not satisfactory and we need to change it. For example, weights can be based on known relative accuracy of observations (when data are averages of samples of different size, the sample sizes are often taken as weights). The weights must take non-negative values. Zero weights are possible, but somewhat extreme (they exclude the given measurement from the analysis), and we do not recommend using them (selected observations can be excluded from the analysis in a much more elegant manner).

1.3 Conventions

The book uses several font types. We use them for distinguishing the basic text of the book from software commands (and other key words) of the R language. For the names of commands and their arguments, we use the bold **Courier New** font. The names of objects that we create during analyses are typed in the normal Courier New font. The other text is typed in the Times New Roman font. Names of variables, parameter values and mathematical formulas are written in *italics*, while the names of factor levels are enclosed in quotation marks. The names of packages are underlined.

For transcribing everything that takes place in the command window of the R environment, we use the Courier New font of a smaller size. For better orientation, we differentiate between commands entered by the user, which we write in bold (for example, **a<-1:5**), and program responses in normal style (for example, a). To save space, some rows of output were cut off and substituted with three dots.

Plots made at the beginning of analyses were created using as few commands as possible and that is why they often do not have appropriate captions, legends, etc. Only plots made at the end of the analysis are closer to real presentation-quality and hence include various details (at the cost of longer commands).

The use of the natural logarithm (with base e) is prevalent in the statistics. We will thus label it simply by log. You need to be aware of the fact that this symbol can have a different meaning elsewhere (for example, in MS Excel it means the decadic logarithm – with base 10).

There are many commercial as well as non-commercial programs available for data analyses. They widely differ by their quality, extent and price. Many inexperienced users work just with the basic software, for example, MS Excel. Some of the popular specialised programs include Statistica or JMP. And finally, there are large (and expensive) packages, such as SPSS or SAS. The difference among statistical software packages is not only in the number of implemented analysis types, but also in their flexibility, i.e. in the possibilities of user programming and other “customised” modifications and the automation of selected procedures. A simple analysis can, of course, be executed using any of the statistical software, however, more advanced methods are offered only by specialised packages.

2.1 The R environment

The good news is that one of the best statistical software packages for modern data analyses is available completely for free. It is called R (R Core Team 2015) and this book is based on its extensive possibilities. It actually covers a much wider area than is presented in this book.

The R programming environment is similar to the commercial program S-Plus (© Insightful Corporation), which used the S programming language; it was developed in the 1980’s in the American AT&T Bell Laboratories (it has been improved several times since then). R uses the dialect of the S language in combination with the Scheme language, which means that, from the user’s point of view, the overwhelming majority of the language for R and S-Plus are either the same or similar. R was created by New Zealanders Ihaka & Gentleman (1996). Today, it is managed by a group of people from around the globe, who call themselves the R Core Team. Thanks to that, the development of R is extremely dynamic. It is constantly evolving and self-improving.

On its own, R is not a user-friendly program of the MS Excel or Statistica type. You will not find nice, colourful windows, pull-down menus and clickable buttons, basically an environment where the main control instrument is the mouse. Be prepared that when you open R, you will only see an empty grey window (Fig. 2-1). R is mainly controlled by commands entered on a keyboard.

There are specialised environments available which can provide a standard user higher comfort when using R (for example, the Rstudio downloadable from <http://www.rstudio.com/>). As we are concentrating on the GLM models and data analysis here, we will not use them in this book.

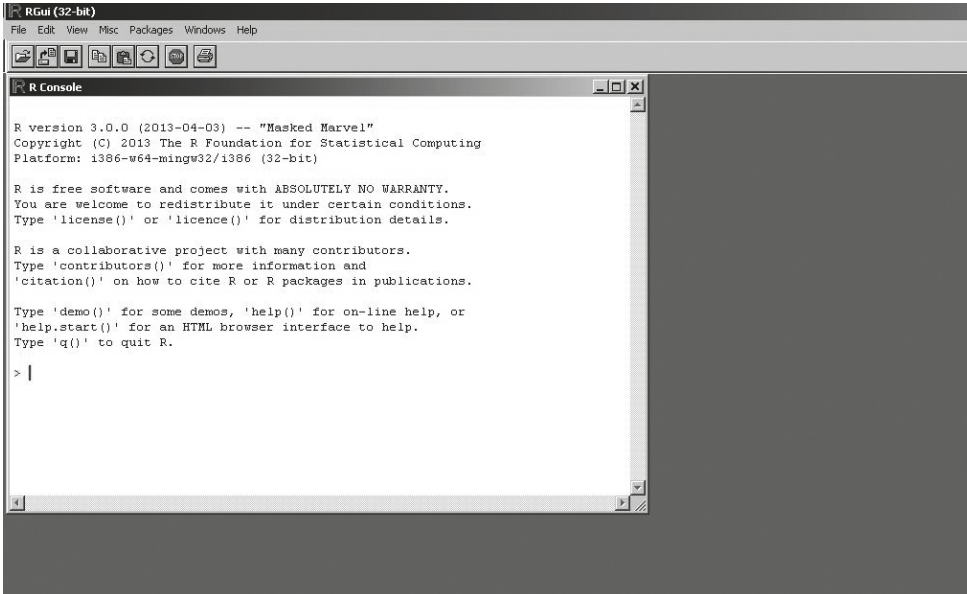


Fig. 2-1 Window of the R environment.

You may ask why we have chosen such “unfriendly” software. We have several reasons for it:

- R is one of the most extensive statistical packages that contain all essential types for modern analyses, which are continuously amended thanks to the continuous work of many people from around the world. Commercial statistical software is usually a closed system (it can be expanded only by purchasing yet another version), while R has been built as an open system. It means that new and additional methods can be added easily and for free at any time.
- Control of this program is not (for a non-statistician) that easy. However, that can be, paradoxically, an advantage because it forces its users to acquire certain knowledge about what they are doing and to think about it (at least) a little bit. Data analyses using commercial software can be somewhat “dangerous” since overly-friendly software allows people who have no idea what they are doing to perform an analysis by a sequence of more or less random clicks. Needless to say that the results of such analyses are often misleading or completely wrong.
- Yet another reason is that “friendly” statistical software will produce a huge volume of information, which they will spout out to you once the analysis is completed. It may not be easy for you to make sense of it. The philosophy of R is completely different – it will display only what the user asks for with the commands. This approach is based on the assumption that users will only use commands that they know something about or, alternatively, that they will look for in help pages or reference literature, thus being stimulated to study the given topic. By default, R outputs are typically quite modest in size. And, in fact, that