



Shravan Vasishth · Michael Broe

The Foundations of Statistics: A Simulation- based Approach

 Springer

The Foundations of Statistics: A Simulation-based Approach

Shravan Vasishth · Michael Broe

The Foundations of Statistics: A Simulation-based Approach

 Springer

Shravan Vasishth
Department of Linguistics
University of Potsdam
Karl-Liebknecht-Str. 24-25
14476 Potsdam
Germany
vasishth@uni-potsdam.de

Michael Broe
Department of Evolution,
Ecology & Organismal Biology
Ohio State University
1304 Museum of Biological Diversity
Kinnear Road 1315
OH 43212 Columbus
USA
broe.1@osu.edu

ISBN 978-3-642-16312-8 e-ISBN 978-3-642-16313-5
DOI 10.1007/978-3-642-16313-5
Springer Heidelberg Dordrecht London New York

© Springer-Verlag Berlin Heidelberg 2011

This work is subject to copyright. All rights are reserved, whether the whole or part of the material is concerned, specifically the rights of translation, reprinting, reuse of illustrations, recitation, broadcasting, reproduction on microfilm or in any other way, and storage in data banks. Duplication of this publication or parts thereof is permitted only under the provisions of the German Copyright Law of September 9, 1965, in its current version, and permission for use must always be obtained from Springer. Violations are liable to prosecution under the German Copyright Law.

The use of general descriptive names, registered names, trademarks, etc. in this publication does not imply, even in the absence of a specific statement, that such names are exempt from the relevant protective laws and regulations and therefore free for general use.

Cover design: WMXDesign GmbH

Cover image: Daniel A. Becker

Printed on acid-free paper

Springer is part of Springer Science+Business Media (www.springer.com)

*SV dedicates this book to his son, Atri;
MB dedicates this book to his parents.*

Foreword

Teaching the fundamental concepts and ideas of statistics is a relatively easy task when students have already completed courses in probability theory and calculus. Such students not only come well-prepared to an introductory statistics course, they are likely to follow additional, more advanced courses in mathematics and related areas. If so, many of the statistical techniques acquired remain relevant and are developed by further training.

For students without much of a background in probability or calculus, the situation is quite different. For many of us, just about everything is new, from integration to probability density functions, and from slopes of regression lines to random variables. Traditionally, introductory statistics courses targeting students from the arts and social sciences seek to explain the basic concepts and results of mathematical statistics, while omitting proofs. To survive such a course, one typically dutifully memorizes the equation for the density of the normal curve, the definition of the central limit theorem, and the sums of squares of a one-way analysis of variance. Once through this bottleneck, there is the safe haven of menu-driven software packages that will calculate everything one needs (and more). Sadly enough, many students will not come to appreciate the beauty of the fundamentals of statistics, and they will also remain somewhat insecure about what the basic concepts actually mean.

The approach taken by Shraavan Vasishth and Michael Broe provides a much more interesting, exciting, and I believe lasting learning experience for students in the social sciences and the arts. (In fact, I believe that students in the sciences will also enjoy lots of the R code in this book). Simulation is a wonderful way of demystifying concepts and techniques that would otherwise remain abstract and distant. By playing around with simulated data, the reader can quickly get an intuitive feel for basic concepts such as the sampling distribution of the sample mean. With the tools provided in this book, the reader can indeed begin to explore the foundations of statistics, and will discover that statistics is actually fun and rewarding. Along the way, the reader will also acquire some basic programming skills, and will be well prepared to use the R software environment for statistical computing and

graphics. R is attractive not only because it is open source, and comes with thousands of add-on packages for specialized data analysis. R is also a great choice because of the many excellent books introducing statistical techniques in a reader-friendly and accessible way, with lots of example data sets. With this introduction, Shravan Vasishth and Michael Broe complement the existing R literature with a very useful and enjoyable hands-on introduction to statistics, doing analysis by synthesis.

Edmonton, October 2010

R. Harald Baayen

Preface

Statistics and hypothesis testing are routinely used in areas that are traditionally not mathematically demanding (an example is psycholinguistics). In such fields, when faced with experimental data in any form, many students and researchers tend to rely on commercial packages to carry out statistical data analysis, often without acquiring much understanding of the logic of statistics they rely on. There are two major problems with this approach. First, the results are often misinterpreted. Second, users are rarely able to flexibly apply techniques relevant to their own research – they use whatever they happened to have learnt from their advisors, and if a slightly new data analysis situation arises, they are unable to use a different method.

A simple solution to the first problem is to teach the foundational ideas of statistical hypothesis testing without using too much mathematics. In order to achieve this, statistics instructors routinely present simulations to students in order to help them intuitively understand things like the Central Limit Theorem. This approach appears to facilitate understanding, but this understanding is fleeting. A deeper and more permanent appreciation of the foundational ideas can be achieved if students re-run and modify the simulations themselves outside the class.

This book is an attempt to address the problem of superficial understanding. It provides a largely non-mathematical, simulation-based introduction to basic statistical concepts, and encourages the reader to try out the simulations themselves using the code provided on the course homepage <http://www.purl.oclc.org/NET/vasishth/VB/>. Since the exercises provided in the text almost always require the use of programming constructs previously introduced, the diligent student acquires basic programming ability as a side effect. This helps to build up the confidence necessary for carrying out more sophisticated analyses. The present book can be considered as the background material necessary for more advanced courses in statistics.

The vehicle for simulation is a freely available software package, R (see the [CRAN website](#) for further details). This book is written using `Sweave`

(pronounced S-weave), which was developed by Leisch, 2002. This means that L^AT_EX and R code are interwoven together.

The style of presentation used in this book is based on a course developed by Michael Broe in the Linguistics department of The Ohio State University. The first author (SV) was a student at the time and attended Michael's course in 2000; later, SV extended the book in the spirit of the original course (which was prepared using commercially available software). Both authors collaborated on the final text.

SV has used this book to teach linguistics undergraduate and graduate students at the University of Saarland, the University of Potsdam, and at the European Summer Schools for Language, Logic and Information held in Edinburgh (2005) and Bordeaux (2009). These courses have shown that the highly motivated student with little to no programming ability and/or mathematical/statistical training can understand everything presented here, and can move on to using R and statistics productively and sensibly.

The book is designed for self-instruction or to accompany a statistics course that involves the use of computers. Some of the examples are from linguistics, but this does not affect the content, which is of general relevance to any scientific discipline. The reader will benefit, as we did, by working through the present book while also consulting some of the books we relied on, in particular Rietveld & van Hout, 2005; Maxwell & Delaney, 2000; Baayen, 2008; Gelman & Hill, 2007.

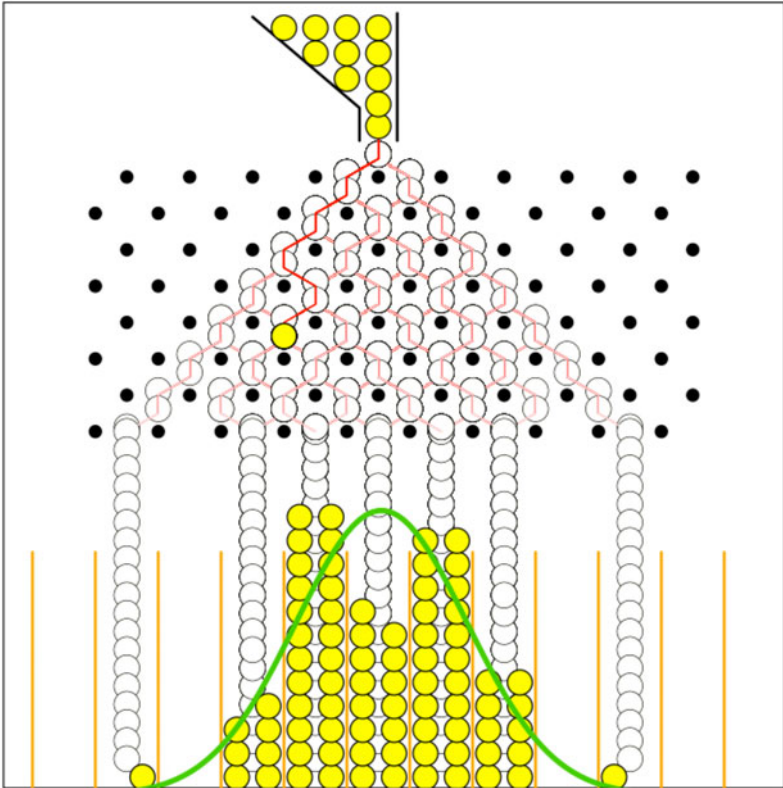
We do not aspire to teach R per se in this book; if this book is used for self-instruction, the reader is expected to either take the initiative themselves to acquire a basic understanding of R, and if this book is used in a taught course, the first few lectures should be devoted to a simple introduction to R.

After completing this book, the reader will be ready to use more advanced books like Gelman and Hill's *Data analysis using regression and multilevel/hierarchical models*, Baayen's *Analyzing Linguistic Data*, and the [online lecture notes by Roger Levy](#).

A lot of people were directly or indirectly involved in the creation of this book. Thanks go to Olga Chiarcos and Federica Corradi dell'Acqua at Springer for patiently putting up with delays in the preparation of this book. In particular, without Olga's initiative and efforts, this book would not have appeared in printed form. SV thanks Reinhold Kliegl, Professor of Psychology at the University of Potsdam, for generously sharing his insights into statistical theory and linear mixed models in particular, and for the opportunity to co-teach courses on statistical data analysis with him; his comments also significantly improved chapter 7. Harald Baayen carefully read the entire book and made important suggestions for improvement; our thanks to him for taking the time. Thanks also to the many students (among them: Pavel Logačev, Rukshin Shaher, and Titus von der Malsburg) who commented on earlier versions of this book. SV is grateful to Andrea Vasishth for support in every aspect of life.

MB thanks the students at The Ohio State University who participated in the development of the original course, and Mary Tait for continued faith and support.

We end with a note on the book cover; it shows a visualization of Galton's box (design by Daniel A. Becker). It is also known as the bean machine or quincunx. This was originally a mechanical device invented by Sir Francis Galton to demonstrate the normal distribution and the central limit theorem. The reader can play with the quincunx using the R version (written by Andrej Blejec) of the simulation, shown below (the code is available from Dr. Blejec's homepage and from the source code accompanying this book). The result of Dr. Blejec's code is shown below.



Berlin, Germany; and Columbus, OH
October 2010

Shravan Vasishth
Michael Broe

Contents

1	Getting Started	1
1.1	Installation: R, L ^A T _E X, and Emacs	1
1.2	How to read this book	2
1.3	Some Simple Commands in R	2
1.4	Graphical Summaries	6
2	Randomness and Probability	9
2.1	Elementary Probability Theory.....	9
2.1.1	The Sum and Product Rules	9
2.1.2	Stones and Rain: A Variant on the Coin-toss Problem.....	11
2.2	The Binomial Distribution	19
2.3	Balls in a Box	22
2.4	Standard Deviation and Sample Size	33
2.4.1	Another Insight: Mean Minimizes Variance	36
2.5	The Binomial versus the Normal Distribution	39
	Problems	41
3	The Sampling Distribution of the Sample Mean	43
3.1	The Central Limit Theorem	47
3.2	σ and $\sigma_{\bar{x}}$	49
3.3	The 95% Confidence Interval for the Sample Mean.....	50
3.4	Realistic Statistical Inference	52
3.5	s is an Unbiased Estimator of σ	52
3.6	The t-distribution	55
3.7	The One-sample t-test	56
3.8	Some Observations on Confidence Intervals	57
3.9	Sample SD, Degrees of Freedom, Unbiased Estimators.....	61
3.10	Summary of the Sampling Process	63
3.11	Significance Tests.....	64
3.12	The Null Hypothesis	65
3.13	z-scores	66

3.14	P-values	67
3.15	Hypothesis Testing: A More Realistic Scenario	71
3.16	Comparing Two Samples	75
3.16.1	H_0 in Two-sample Problems	76
	Problems	79
4	Power	81
4.1	Hypothesis Testing Revisited	81
4.2	Type I and Type II Errors	82
4.3	Equivalence Testing	91
4.3.1	Equivalence Testing Example	91
4.3.2	TOST Approach to the Stegner et al. Example	92
4.3.3	Equivalence Testing Example: CIs Approach	94
4.4	Observed Power and Null Results	94
	Problems	96
5	Analysis of Variance (ANOVA)	97
5.1	Comparing Three Populations	97
5.2	ANOVA	99
5.2.1	Statistical Models	100
5.2.2	Variance of Sample Means as a Possible Statistic	103
5.2.3	Analyzing the Variance	105
5.3	Hypothesis Testing	111
5.3.1	MS-within, MS-between as Statistics	112
5.3.2	The F-distribution	113
5.3.3	ANOVA in R	118
5.3.4	MS-within, Three Non-identical Populations	118
5.3.5	The F-distribution with Unequal Variances	120
5.4	ANOVA as a Linear Model	121
	Problems	125
6	Bivariate Statistics and Linear Models	127
6.1	Variance and Hypothesis Testing for Regression	137
6.1.1	Sum of Squares and Correlation	142
	Problems	142
7	An Introduction to Linear Mixed Models	145
7.1	Introduction	145
7.2	Simple Linear Model	146
7.3	The Levels of the Complex Linear Model	155
7.4	Further Reading	159

A	Random Variables	161
A.1	The Probability Distribution in Statistical Inference	162
A.2	Expectation	162
A.3	Properties of Expectation	164
A.4	Variance	165
A.5	Important Properties of Variance	165
A.6	Mean and SD of the Binomial Distribution	166
A.7	Sample versus Population Means and Variances	167
A.8	Summing up	168
	Problems	169
B	Basic R Commands and Data Structures	171
	References	175
	Index	177

Chapter 1

Getting Started

The main goal of this book is to help you understand the principles behind inferential statistics, and to use and customize statistical tests to your needs. The vehicle for this will be a programming language called R (google for ‘CRAN,’ which stands for the Comprehensive R Archive Network). Let’s start with installing R and related software.

1.1 Installation: R, L^AT_EX, and Emacs

In order to use the code that comes with this book, you only need to install R. The latest version can be downloaded from the [CRAN website](#). However, other freely available software provides a set of tools that work together with R to give a very pleasant computing environment. The least that you need to know about is L^AT_EX, Emacs, and Emacs Speaks Statistics. Other tools that will further enhance your working experience with L^AT_EX are AucTeX, RefTeX, preview-latex, and Python. None of these are required but are highly recommended for typesetting and other sub-tasks necessary for data analysis.

There are many advantages to using R with these tools. For example, R and L^AT_EX code can be intermixed in Emacs using noweb mode. R can output data tables etc. in L^AT_EX format, allowing you to efficiently integrate your scientific writing with the data analysis. This book was typeset using many of the above tools. For more on L^AT_EX, see the <http://www.tug.org>, the TeX Users Group website.

The installation of this working environment differs from one operating system to another. In Linux-like environments, most of these tools are already pre-installed. For Windows and Macintosh you will need to read the manual pages on the CRAN website.

After you have installed R on your machine, the second thing you need to do before proceeding any further with this book is to learn a little bit about R. The present book is not intended to be an introduction to R. For short,

comprehensive and freely available introductions, look at the Manuals on the R homepage, and particularly under the link ‘Contributed.’ You should spend a few hours or even days studying some of the shorter articles in the Contributed section of the CRAN archive. In particular, you need to know basic things like starting up R, simple arithmetic operations, and quitting R. It is possible to skip this step and to learn R as you read this book, but in that case you have to be prepared to look up the online help available with R. For readers who want to start using the present book immediately, we provide a very basic introduction to R in Appendix B; the reader should work through this material before reading further.

1.2 How to read this book

We recommend that the book be read with an R session open on your computer and the code accompanying this book be kept open as an R file. The code used in this book is available from the homepage:

<http://www.purl.oclc.org/NET/vasishth/VB/>

You will get the most out of this book if you run the code as you read along, pausing to experiment with the code (changing parameters, asking yourself: “what would happen if I changed this setting?”, etc.). Passive reading of the textbook will probably not yield much. Do not hesitate to re-read chapters! This material is best digested by revisiting it several times. The book chapters are intended to be read in order, since the later chapters build on concepts introduced in earlier parts of the book.

The accompanying website for the book contains (among other things):

1. A blog for asking questions that the authors or other readers can answer, or for submitting comments (corrections, suggestions, etc.) regarding the book that other readers can also benefit from.
2. Links to useful websites and other material relating to statistics.
3. Additional material that contains a more advanced treatment of some of the issues discussed in the textbook.

1.3 Some Simple Commands in R

We begin with a short session that aims to familiarize you with R and very basic interaction with data.

Let’s assume for argument’s sake that we have the grades of eleven students in a final examination for a statistics course. Both the instructor and the students are probably interested in finding out at least the maximum and

minimum scores. But hidden in these scores is much more information about the students' grades.

Assuming a maximum possible score of 100, let's first start up R and input the scores (which we just made up).

```
> scores <- c(99, 97, 72, 56, 88, 80, 74, 95, 66,
              57, 89)

[1] 99 97 72 56 88 80 74 95 66 57 89
```

(When you type the above on the R command line you will not see R echo the contents of `scores`. You can type `scores` on the command line to print out the contents. Another way to print out the content of a command is to wrap it in parentheses: `(scores)`. We are able to print out the contents of `scores` above without using either of these methods because we are using Sweave with L^AT_EX, and Sweave allows the user to print when needed. See the Sweave homepage, <http://www.stat.uni-muenchen.de/~leisch/Sweave/>, for details)

Now we ask the following questions using R: (a) what's the maximum score? (b) what's the minimum?

```
> max(scores)

[1] 99

> min(scores)

[1] 56
```

We could stop here. But there is much more information in this simple dataset, and it tells us a great deal more about the students than the maximum and minimum grades.

The first thing we can ask is: what is the average or mean score? For any collection of numbers, their MEAN is the sum of the numbers divided by the length of the vector:

$$\bar{x} = \frac{x_1 + x_2 + \cdots + x_n}{n} = \frac{1}{n} \sum_{i=1}^n x_i \quad (1.1)$$

The notation $\sum_{i=1}^n$ is simply an abbreviation for the statement that the numbers going from x_1 to x_n should be added up.

The mean tells you something interesting about that collection of students: if they had all scored high marks, say in the 90's, the mean would be high, and if not then it would be relatively low. The mean gives you one number that summarizes the data succinctly. We can ask R to compute the mean as follows:

```
> mean(scores)
```