Monographs on Statistics and Applied Probability 57

An Introduction to the Bootstrap

Bradley Efron Robert J. Tibshirani

CHAPMAN & HALL/CRC

MONOGRAPHS ON STATISTICS AND APPLIED PROBABILITY

General Editors

D.R. Cox, D.V. Hinkley, N. Reid, D.B. Rubin and B.W. Silverman

1 Stochastic Population Models in Ecology and Epidemiology M.S. Bartlett (1960)

2 Queues D.R. Cox and W.L. Smith (1961)

3 Monte Carlo Methods J.M. Hammersley and D.C. Handscomb (1964)

4 The Statistical Analysis of Series of Events D.R. Cox and P.A.W. Lewis (1966)

5 Population Genetics W.J. Ewens (1969)

6 Probability, Statistics and Time M.S. Bartlett (1975)

7 Statistical Inference S.D. Silvey (1975)

8 The Analysis of Contingency Tables B.S. Everitt (1977)

9 Multivariate Analysis in Behavioural Research A.E. Maxwell (1977)

10 Stochastic Abundance Models S. Engen (1978)

11 Some Basic Theory for Statistical Inference EJ.G. Pitman (1979)

12 Point Processes D.R. Cox and V. Isham (1980)

13 Identification of Outliers D.M. Hawkins (1980)

14 Optimal Design S.D. Silvey (1980)

15 Finite Mixture Distributions B.S. Everitt and D.J. Hand (1981)

16 Classification A.D. Gordon (1981)

17 Distribution-free Statistical Methods J.S. Mariz (1981)

18 Residuals and Influence in Regression R.D. Cook and S. Weisberg (1982)

19 Applications of Queueing Theory G.F. Newell (1982)

20 Risk Theory, 3rd edition R.E. Beard, T. Pentikainen and E. Pesonen (1984)

21 Analysis of Survival Data D.R. Cox and D. Oakes (1984)

22 An Introduction to Latent Variable Models B.S. Everitt (1984)

23 Bandit Problems D.A. Berry and B. Fristedt (1985)

24 Stochastic Modelling and Control M.H.A. Davis and R. Vinter (1985)

25 The Statistical Analysis of Compositional Data J. Aitchison (1986)

26 Density Estimation for Statistical and Data Analysis B.W. Silverman (1986)

27 Regression Analysis with Applications B.G. Wetherill (1986)

28 Sequential Methods in Statistics, 3rd edition G.B. Wetherill (1986)
29 Tensor methods in Statistics P. McCullagh (1987)

30 Transformation and Weighting in Regression R.J. Carroll and D. Ruppert (1988)

31 Asymptotic Techniques for Use in Statistics O.E. Barndoff-Nielson and D.R. Cox (1989)

Analysis of Binary Data, 2nd edition D.R. Cox and E.J. Snell (1989)
 Analysis of Infectious Disease Data N.G. Becker (1989)

- 34 Design and Analysis of Cross-Over Trials B. Jones and M.G. Kenward (1989)
 35 Empirical Bayes Method, 2nd edition J.S. Maritz and T. Lwin (1989)
 - 36 Symmetric Multivariate and Related Distributions K.-T. Fang, S. Kotz and K. Ng (1989)
- 37 Generalized Linear Models, 2nd edition P. McCullagh and J.A. Nelder (1989)
 38 Cyclic Designs J.A. John (1987)
 - 39 Analog Estimation Methods in Econometrics C.F. Manski (1988)
 40 Subset Selection in Regression A.J. Miller (1990)
 - 41 Analysis of Repeated Measures M. Crowder and D.J. Hand (1990)
 - 42 Statistical Reasoning with Imprecise Probabilities P. Walley (1990)
 - 43 Generalized Additive Models T.J. Hastie and R.J. Tibshirani (1990)
- 44 Inspection Errors for Attributes in Quality Control N.L. Johnson, S. Kotz and X. Wu (1991)
 - 45 The Analysis of Contingency Tables, 2nd edition B.S. Everitt (1992)
 46 The Analysis of Quantal Response Data B.J.T. Morgan (1992)
 - 47 Longitudinal Data with Serial Correlation: A State-Space Approach R.H. Jones (1993)
 - 48 Differential Geometry and Statistics M.K. Murray and J.W. Rice (1993)
 49 Markov Models and Optimization M.H.A. Davies (1993)
 - 50 Chaos and Networks: Statistical and Probabilistic Aspects Edited by O. Barndorff-Nielsen et al. (1993)
 - 51 Number Theoretic Methods in Statistics K.-T. Fang and W. Yuan (1993)
 - 52 Inference and Asymptotics O. Barndorff-Nielsen and D.R. Cox (1993)
 - 53 Practical Risk Theory for Actuaries C.D. Daykin, T. Pentikainen and M. Pesonen (1993)
 - 54 Statistical Concepts and Applications in Medicine J. Aitchison and IJ. Lauder (1994)

55 Predictive Inference S. Geisser (1993)

- 56 Model-Free Curve Estimation M. Tarter and M. Lock (1993)
- 57 An Introduction to the Bootstrap B. Efron and R. Tibshirani (1993)
- (Full details concerning this series are available from the Publishers.)



Bradley Efron

Department of Statistics Stanford University and

Robert J. Tibshirani

Department of Preventative Medicine and Biostatistics and Department of Statistics, University of Toronto

CHAPMAN & HALL/CRC

Boca Raton London New York Washington, D.C.

Chapman & Hall/CRC Taylor & Francis Group 6000 Broken Sound Parkway NW, Suite 300 Boca Raton, FL 33487-2742 Milton Park, Abingdon Oxon OX 14 4RN

© 1994 by Taylor & Francis Group, LLC Chapman & Hall/CRC is an imprint of Taylor & Francis Group

No claim to original U.S. Government works Printed in the United States of America on acid-free paper 25 24 23 22 21 20 19 18 17 16 15 14 13 International Standard Book Number-13: 978-0-412-04231-7 (Hardcover)

This book contains information obtained from authentic and highly regarded sources. Reasonable efforts have been made to publish reliable data and information, but the author and publisher cannot assume responsibility for the validity of all materials or the consequences of their use. The authors and publishers have attempted to trace the copyright holders of all material reproduced in this publication and apologize to copyright holders if permission to publish in this form has not been obtained. If any copyright material has not been acknowledged please write and let us know so we may rectify it in any future reprint.

No part of this book may be reprinted, reproduced, transmitted, or utilized in any form by any electronic, mechanical, or other means, now known or hereafter invented, including photocopying, microfilming, and recording, or in any information storage or retrieval system, without written permission from the publishers.

For permission to photocopy or use material electronically from this work, please access www.copyright.com (http://www.copyright.com/) or contact the Copyright Clearance Center, Inc. (CCC) 222 Rosewood Drive, Danvers, MA 01923, 978-750-8400. CCC is a not-for-profit organization that provides licenses and registration for a variety of users. For organizations that have been granted a photocopy license by the CCC, a separate system of payment has been arranged.

Trademark Notice: Product or corporate names may be trademarks or registered trademarks, and are used only for identification and explanation without intent to infringe.

Library of Congress Cataloging-in-Publication Data

93-4489

Efron, Bradley. An introduction to the bootstrap/Brad Efron, Rob Tibshirani. p. cm. Includes bibliographical references and index. ISBN 0-412-04231-2 1. Bootstrap (Statistics). I. Tibshirani, Robert. II. Title. QA276.8.E3745 1993 519.5'44—dc20

Visit the Taylor & Francis Web site at http://www.taylorandfrancis.com

and the CRC Press Web site at http://www.crcpress.com

TO CHERYL, CHARLIE, RYAN AND JULIE

> AND TO THE MEMORY OF RUPERT G. MILLER, JR.



Contents

P	Preface					
1	Introduction					
	1.1	An overview of this book	6			
	1.2	Information for instructors	8			
	1.3	Some of the notation used in the book	9			
2	The	e accuracy of a sample mean	10			
	2.1	Problems	15			
3	·Raı	ndom samples and probabilities	17			
	3.1	Introduction	17			
	3.2	Random samples	17			
	3.3	Probability theory	20			
	3.4	Problems	28			
4	The	e empirical distribution function and the plug-	in			
	priı	nciple	31			
	4.1	Introduction	31			
	4.2	The empirical distribution function	31			
	4.3	The plug-in principle	35			
	4.4	Problems	37			
5	Sta	ndard errors and estimated standard errors	39			
	5.1	Introduction	39			
	5.2	The standard error of a mean	39			
	5.3	Estimating the standard error of the mean	42			
	5.4	Problems	43			

6	The	bootstrap estimate of standard error	45
	6.1	Introduction	45
	6.2	The bootstrap estimate of standard error	45
	6.3	Example: the correlation coefficient	49
	6.4	The number of bootstrap replications B	50
	6.5	The parametric bootstrap	53
	6.6	Bibliographic notes	56
	6.7	Problems	57
7	Boo	tstrap standard errors: some examples	60
	7.1	Introduction	60
	7.2	Example 1: test score data	61
	7.3	Example 2: curve fitting	70
	7.4	An example of bootstrap failure	81
	7.5	Bibliographic notes	81
	7.6	Problems	82
8	Mor	e complicated data structures	86
	8.1	Introduction	86
	8.2	One-sample problems	86
	8.3	The two-sample problem	88
	8.4	More general data structures	90
	8.5	Example: lutenizing hormone	92
	8.6	The moving blocks bootstrap	99
	8.7	Bibliographic notes	102
	8.8	Problems	103
9	Reg	ression models	105
	9.1	Introduction	105
	9.2	The linear regression model	105
	9.3	Example: the hormone data	107
	9.4	Application of the bootstrap	111
	9.5	Bootstrapping pairs vs bootstrapping residuals	113
	9.6	Example: the cell survival data	115
	9.7	Least median of squares	117
	9.8	Bibliographic notes	121
	9.9	Problems	121
10	Esti	mates of bias	124
	10.1	Introduction	124

viii

CONTENTS	ix
10.2 The bootstrap estimate of bias	124
10.3 Example: the patch data	126
10.4 An improved estimate of bias	130
10.5 The jackknife estimate of bias	133
10.6 Bias correction	138
10.7 Bibliographic notes	139
10.8 Problems	139
11 The jackknife	141
11.1 Introduction	141
11.2 Definition of the jackknife	141
11.3 Example: test score data	143
11.4 Pseudo-values	145
11.5 Relationship between the jackknife and bootstrap	145
11.6 Failure of the jackknife	148
11.7 The delete- d jackknife	149
11.8 Bibliographic notes	149
11.9 Problems	150
12 Confidence intervals based on bootstrap "tables"	153
12.1 Introduction	153
12.2 Some background on confidence intervals	155
12.3 Relation between confidence intervals and hypothe-	
sis tests	156
12.4 Student's t interval	158
12.5 The bootstrap- t interval	160
12.6 Transformations and the bootstrap- t	162
12.7 Bibliographic notes	166
12.8 Problems	166
13 Confidence intervals based on bootstrap	
percentiles	168
13.1 Introduction	168
13.2 Standard normal intervals	168
13.3 The percentile interval	170
13.4 Is the percentile interval backwards?	174
13.5 Coverage performance	174
13.6 The transformation-respecting property	175
13.7 The range-preserving property	176
13.8 Discussion	176

x	CON	TENTS
	13.9 Bibliographic notes	176
	13.10 Problems	177
14	Better bootstrap confidence intervals	178
	14.1 Introduction	178
	14.2 Example: the spatial test data	179
	14.3 The BC_a method	184
	14.4 The ABC method	188
	14.5 Example: the tooth data	190
	14.6 Bibliographic notes	199
	14.7 Problems	199
15	Permutation tests	202
	15.1 Introduction	202
	15.2 The two-sample problem	202
	15.3 Other test statistics	210
	15.4 Relationship of hypothesis tests to confidence	
	intervals and the bootstrap	214
	15.5 Bibliographic notes	218
	15.6 Problems	218
16	Hypothesis testing with the bootstrap	220
	16.1 Introduction	220
	16.2 The two-sample problem	220
	16.3 Relationship between the permutation test and the	e
	bootstrap	223
	16.4 The one-sample problem	224
	16.5 Testing multimodality of a population	227
	16.6 Discussion	232
	16.7 Bibliographic notes	233
	16.8 Problems	234
17	Cross-validation and other estimates of predict	ion
	error	237
	17.1 Introduction	237
	17.2 Example: hormone data	238
	17.3 Cross-validation	239
	17.4 C_p and other estimates of prediction error	242
	17.5 Example: classification trees	243
	17.6 Bootstrap estimates of prediction error	247

CO	NTEN	JTS	xi
		17.6.1 Overview	247
		17.6.2 Some details	249
	17.7	The .632 bootstrap estimator	252
	17.8	Discussion	254
	17.9	Bibliographic notes	255
	17.10) Problems	255
18	Ada	ptive estimation and calibration	258
	18.1	Introduction	258
	18.2	Example: smoothing parameter selection for curve fitting	258
	18.3	Example: calibration of a confidence point	263
	18.4	Some general considerations	266
	18.5	Bibliographic notes	268
	18.6	Problems	269
19	Ass	essing the error in bootstrap estimates	271
	19.1	Introduction	271
	19.2	Standard error estimation	272
	19.3	Percentile estimation	273
	19.4	The jackknife-after-bootstrap	275
	19.5	Derivations	280
	19.6	Bibliographic notes	281
	19.7	Problems	281
20	Ag	eometrical representation for the bootstrap and	d
	jack	knite	283
	20.1	Introduction	283
	20.2	Bootstrap sampling	285
	20.3	The jackknife as an approximation to the bootstrap	287
	20.4	Other jackknife approximations	289
	20.5	Estimates of bias	290
	20.6	An example	293
	20.7	Bibliographic notes	295
	20.8	Problems	295
21	An	overview of nonparametric and parametric	9 0 <i>e</i>
	1111e	Introduction	490 204
	21.1 91 9	Distributions densities and likelihood functions	290
	21.2	Distributions, densities and likelihood functions	290

	21.3 Functional statistics and influence functions	298
	21.4 Parametric maximum likelihood inference	302
	21.5 The parametric bootstrap	306
	21.6 Relation of parametric maximum likelihood, boot-	
	strap and jackknife approaches	307
	21.6.1 Example: influence components for the mean	309
	21.7 The empirical cdf as a maximum likelihood estimate	310
	21.8 The sandwich estimator	310
	21.8.1 Example: Mouse data	311
	21.9 The delta method	313
	21.9.1 Example: delta method for the mean	315
	21.9.2 Example: delta method for the correlation	
	$\operatorname{coefficient}$	315
	21.10 Relationship between the delta method and in-	
	finitesimal jackknife	315
	21.11 Exponential families	316
	21.12 Bibliographic notes	319
	21.13 Problems	320
22	Further topics in bootstrap confidence intervals	321
	22.1 Introduction	321
	22.2 Correctness and accuracy	321
	22.3 Confidence points based on approximate pivots	322
	22.4 The BC_a interval	325
	22.5 The underlying basis for the BC_a interval	326
	22.6 The ABC approximation	328
	22.7 Least favorable families	331
	22.8 The ABC_q method and transformations	333
	22.8 The ABC _q method and transformations 22.9 Discussion	333 334
	 22.8 The ABC_q method and transformations 22.9 Discussion 22.10 Bibliographic notes 	333 334 335
	 22.8 The ABC_q method and transformations 22.9 Discussion 22.10 Bibliographic notes 22.11 Problems 	333 334 335 335
23	 22.8 The ABC_q method and transformations 22.9 Discussion 22.10 Bibliographic notes 22.11 Problems Efficient bootstrap computations	333 334 335 335 338
23	 22.8 The ABC_q method and transformations 22.9 Discussion 22.10 Bibliographic notes 22.11 Problems Efficient bootstrap computations 23.1 Introduction 	333 334 335 335 335 338 338
23	 22.8 The ABC_q method and transformations 22.9 Discussion 22.10 Bibliographic notes 22.11 Problems Efficient bootstrap computations 23.1 Introduction 23.2 Post-sampling adjustments 	333 334 335 335 335 338 338 340
23	 22.8 The ABC_q method and transformations 22.9 Discussion 22.10 Bibliographic notes 22.11 Problems Efficient bootstrap computations 23.1 Introduction 23.2 Post-sampling adjustments 23.3 Application to bootstrap bias estimation 	333 334 335 335 335 338 338 340 342
23	 22.8 The ABC_q method and transformations 22.9 Discussion 22.10 Bibliographic notes 22.11 Problems Efficient bootstrap computations 23.1 Introduction 23.2 Post-sampling adjustments 23.3 Application to bootstrap bias estimation 23.4 Application to bootstrap variance estimation 	333 334 335 335 335 338 340 342 346
23	 22.8 The ABC_q method and transformations 22.9 Discussion 22.10 Bibliographic notes 22.11 Problems Efficient bootstrap computations 23.1 Introduction 23.2 Post-sampling adjustments 23.3 Application to bootstrap bias estimation 23.4 Application to bootstrap variance estimation 23.5 Pre- and post-sampling adjustments 	333 334 335 335 335 338 340 342 346 348
23	 22.8 The ABC_q method and transformations 22.9 Discussion 22.10 Bibliographic notes 22.11 Problems Efficient bootstrap computations 23.1 Introduction 23.2 Post-sampling adjustments 23.3 Application to bootstrap bias estimation 23.4 Application to bootstrap variance estimation 23.5 Pre- and post-sampling adjustments 23.6 Importance sampling for tail probabilities 	333 334 335 335 338 340 342 346 348 348 349
23	 22.8 The ABC_q method and transformations 22.9 Discussion 22.10 Bibliographic notes 22.11 Problems Efficient bootstrap computations 23.1 Introduction 23.2 Post-sampling adjustments 23.3 Application to bootstrap bias estimation 23.4 Application to bootstrap variance estimation 23.5 Pre- and post-sampling adjustments 23.6 Importance sampling for tail probabilities 23.7 Application to bootstrap tail probabilities 	333 334 335 335 335 338 340 342 346 348 349 352

xii

CONTENTS				
23.8 Bibliographic notes	356			
23.9 Problems	357			
24 Approximate likelihoods	358			
24.1 Introduction	358			
24.2 Empirical likelihood	360			
24.3 Approximate pivot methods	362			
24.4 Bootstrap partial likelihood	364			
24.5 Implied likelihood	367			
24.6 Discussion	370			
24.7 Bibliographic notes	371			
24.8 Problems	371			
25 Bootstrap bioequivalence	372			
25.1 Introduction	372			
25.2 A bioequivalence problem	372			
25.3 Bootstrap confidence intervals	374			
25.4 Bootstrap power calculations	379			
25.5 A more careful power calculation	381			
25.6 Fieller's intervals	384			
25.7 Bibliographic notes	389			
25.8 Problems	389			
26 Discussion and further topics	392			
26.1 Discussion	392			
26.2 Some questions about the bootstrap	394			
26.3 References on further topics	396			
Appendix: software for bootstrap computations	398			
Introduction	398			
Some available software	399			
S language functions	399			
References	413			
Author index				
Subject index	430			

Preface

Dear friend, theory is all gray, and the golden tree of life is green. Goethe, from "Faust"

The ability to simplify means to eliminate the unnecessary so that the necessary may speak.

Hans Hoffmann

Statistics is a subject of amazingly many uses and surprisingly few effective practitioners. The traditional road to statistical knowledge is blocked, for most, by a formidable wall of mathematics. Our approach here avoids that wall. The bootstrap is a computerbased method of statistical inference that can answer many real statistical questions without formulas. Our goal in this book is to arm scientists and engineers, as well as statisticians, with computational techniques that they can use to analyze and understand complicated data sets.

The word "understand" is an important one in the previous sentence. This is not a statistical cookbook. We aim to give the reader a good intuitive understanding of statistical inference.

One of the charms of the bootstrap is the direct appreciation it gives of variance, bias, coverage, and other probabilistic phenomena. What does it mean that a confidence interval contains the true value with probability .90? The usual textbook answer appears formidably abstract to most beginning students. Bootstrap confidence intervals are directly constructed from real data sets, using a simple computer algorithm. This doesn't necessarily make it easy to understand confidence intervals, but at least the difficulties are the appropriate conceptual ones, and not mathematical muddles.

PREFACE

Much of the exposition in our book is based on the analysis of real data sets. The mouse data, the stamp data, the tooth data, the hormone data, and other small but genuine examples, are an important part of the presentation. These are especially valuable if the reader can try his own computations on them. Personal computers are sufficient to handle most bootstrap computations for these small data sets.

This book does not give a rigorous technical treatment of the bootstrap, and we concentrate on the ideas rather than their mathematical justification. Many of these ideas are quite sophisticated, however, and this book is not just for beginners. The presentation starts off slowly but builds in both its scope and depth. More mathematically advanced accounts of the bootstrap may be found in papers and books by many researchers that are listed in the Bibliographic notes at the end of the chapters.

We would like to thank Andreas Buja, Anthony Davison, Peter Hall, Trevor Hastie, John Rice, Bernard Silverman, James Stafford and Sami Tibshirani for making very helpful comments and suggestions on the manuscript. We especially thank Timothy Hesterberg and Cliff Lunneborg for the great deal of time and effort that they spent on reading and preparing comments. Thanks to Maria-Luisa Gardner for providing expert advice on the "rules of punctuation." We would also like to thank numerous students at both Stanford University and the University of Toronto for pointing out errors in earlier drafts, and colleagues and staff at our universities for their support. Thanks to Tom Glinos of the University of Toronto for maintaining a healthy computing environment. Karola DeCleve typed much of the first draft of this book, and maintained vigilance against errors during its entire history. All of this was done cheerfully and in a most helpful manner, for which we are truly grateful. Trevor Hastie provided expert "S" and TFX advice, at crucial stages in the project.

We were lucky to have not one but two superb editors working on this project. Bea Schube got us going, before starting her retirement; Bea has done a great deal for the statistics profession and we wish her all the best. John Kimmel carried the ball after Bea left, and did an excellent job. We thank our copy-editor Jim Geronimo for his thorough correction of the manuscript, and take responsibility for any errors that remain.

The first author was supported by the National Institutes of Health and the National Science Foundation. Both groups have supported the development of statistical theory at Stanford, including much of the theory behind this book. The second author would like to thank his wife Cheryl for her understanding and support during this entire project, and his parents for a lifetime of encouragement. He gratefully acknowledges the support of the Natural Sciences and Engineering Research Council of Canada.

Palo Alto and Toronto June 1993 Bradley Efron Robert Tibshirani

xvi

CHAPTER 1

Introduction

Statistics is the science of learning from experience, especially experience that arrives a little bit at a time. The earliest information science was statistics, originating in about 1650. This century has seen statistical techniques become the analytic methods of choice in biomedical science, psychology, education, economics, communications theory, sociology, genetic studies, epidemiology, and other areas. Recently, traditional sciences like geology, physics, and astronomy have begun to make increasing use of statistical methods as they focus on areas that demand informational efficiency, such as the study of rare and exotic particles or extremely distant galaxies.

Most people are not natural-born statisticians. Left to our own devices we are not very good at picking out patterns from a sea of noisy data. To put it another way, we are all too good at picking out non-existent patterns that happen to suit our purposes. Statistical theory attacks the problem from both ends. It provides optimal methods for finding a real signal in a noisy background, and also provides strict checks against the overinterpretation of random patterns.

Statistical theory attempts to answer three basic questions:

- (1) How should I collect my data?
- (2) How should I analyze and summarize the data that I've collected?
- (3) How accurate are my data summaries?

Question 3 constitutes part of the process known as statistical inference. The bootstrap is a recently developed technique for making certain kinds of statistical inferences. It is only recently developed because it requires modern computer power to simplify the often intricate calculations of traditional statistical theory.

The explanations that we will give for the bootstrap, and other

computer-based methods, involve explanations of traditional ideas in statistical inference. The basic ideas of statistics haven't changed, but their implementation has. The modern computer lets us apply these ideas flexibly, quickly, easily, and with a minimum of mathematical assumptions. Our primary purpose in the book is to explain when and why bootstrap methods work, and how they can be applied in a wide variety of real data-analytic situations.

All three basic statistical concepts, data collection, summary and inference, are illustrated in the New York Times excerpt of Figure 1.1. A study was done to see if small aspirin doses would prevent heart attacks in healthy middle-aged men. The data for the aspirin study were collected in a particularly efficient way: by a controlled, randomized, double-blind study. One half of the subjects received aspirin and the other half received a control substance, or placebo, with no active ingredients. The subjects were randomly assigned to the aspirin or placebo groups. Both the subjects and the supervising physicians were blinded to the assignments, with the statisticians keeping a secret code of who received which substance. Scientists, like everyone else, want the project they are working on to succeed. The elaborate precautions of a controlled, randomized, blinded experiment guard against seeing benefits that don't exist, while maximizing the chance of detecting a genuine positive effect.

The summary statistics in the newspaper article are very simple:

	heart attacks	subjects
	(fatal plus non-fatal)	
aspirin group:	104	11037
placebo group:	189	11034

We will see examples of much more complicated summaries in later chapters. One advantage of using a good experimental design is a simplification of its results. What strikes the eye here is the lower rate of heart attacks in the aspirin group. The ratio of the two rates is

$$\widehat{\theta} = \frac{104/11037}{189/11034} = .55. \tag{1.1}$$

If this study can be believed, and its solid design makes it very believable, the aspirin-takers only have 55% as many heart attacks as placebo-takers.

Of course we are not really interested in $\hat{\theta}$, the estimated ratio. What we would like to know is θ , the true ratio, that is the ratio

HEART ATTACK RISK Found to be cut by taking aspirin

LIFESAVING EFFECTS SEEN

Study Finds Benefit of Tablet Every Other. Day Is Much Greater Than Expected

By HAROLD M. SCHMECK Jr.

A major nationwide study shows that a single aspirin tablet every other day can sharply reduce a man's risk of heart attack and death from heart attack.

The lifesaving effects were so dramatic that the study was halted in mid-December so that the results could be reported as soon as possible to the participants and to the medical profession in general.

The magnitude of the beneficial effect was far greater than expected, Dr. Charles H. Hennekens of Harvard, principal investigator in the research, and in a telephone interview. The risk of myocardial infarction, the technical name for heart attack, was cut almost in half.

'Extreme Beneficial Effect'

A special report said the results showed "a statistically extreme beneficial effect" from the use of aspirin. The report is to be published Thursday in The New England Journal of Medicine.

In recent years smaller studies have demonstrated that a person who has had one heart attack can reduce the risk of a second by taking aspirin, but there had been no proof that the beneficial effect would extend to the general male population.

Dr. Claude Lenfant, the director of the National Heart Lung and Blood institute, said the findings were "extremely important," but he said the general public should not take the report as an indication that everyone should start taking aspirin.

Figure 1.1. Front-page news from the New York Times of January 27, 1987. Reproduced by permission of the New York Times.

we would see if we could treat all subjects, and not just a sample of them. The value $\hat{\theta} = .55$ is only an estimate of θ . The sample seems large here, 22071 subjects in all, but the conclusion that aspirin works is really based on a smaller number, the 293 observed heart attacks. How do we know that $\hat{\theta}$ might not come out much less favorably if the experiment were run again?

This is where statistical inference comes in. Statistical theory allows us to make the following inference: the true value of θ lies in the interval

$$.43 < \theta < .70 \tag{1.2}$$

with 95% confidence. Statement (1.2) is a classical confidence interval, of the type discussed in Chapters 12–14, and 22. It says that if we ran a much bigger experiment, with millions of subjects, the ratio of rates probably wouldn't be too much different than (1.1). We almost certainly wouldn't decide that θ exceeded 1, that is that aspirin was actually harmful. It is really rather amazing that the same data that give us an estimated value, $\hat{\theta} = .55$ in this case, also can give us a good idea of the estimate's accuracy.

Statistical inference is serious business. A lot can ride on the decision of whether or not an observed effect is real. The aspirin study tracked strokes as well as heart attacks, with the following results:

	$\operatorname{strokes}$	subjects	
aspirin group:	119	11037	
placebo group:	98	11034	(1.3)

For strokes, the ratio of rates is

$$\widehat{\theta} = \frac{119/11037}{98/11034} = 1.21. \tag{1.4}$$

It now looks like taking aspirin is actually harmful. However the interval for the true stroke ratio θ turns out to be

$$.93 < \theta < 1.59$$
 (1.5)

with 95% confidence. This includes the neutral value $\theta = 1$, at which aspirin would be no better or worse than placebo vis-à-vis strokes. In the language of statistical hypothesis testing, aspirin was found to be significantly beneficial for preventing heart attacks, but not significantly harmful for causing strokes. The opposite conclusion had been reached in an older, smaller study concerning men

who had experienced previous heart attacks. The aspirin treatment remains mildly controversial for such patients.

The bootstrap is a data-based simulation method for statistical inference, which can be used to produce inferences like (1.2) and (1.5). The use of the term bootstrap derives from the phrase to pull oneself up by one's bootstrap, widely thought to be based on one of the eighteenth century Adventures of Baron Munchausen, by Rudolph Erich Raspe. (The Baron had fallen to the bottom of a deep lake. Just when it looked like all was lost, he thought to pick himself up by his own bootstraps.) It is not the same as the term "bootstrap" used in computer science meaning to "boot" a computer from a set of core instructions, though the derivation is similar.

Here is how the bootstrap works in the stroke example. We create two populations: the first consisting of 119 ones and 11037-119=10918 zeroes, and the second consisting of 98 ones and 11034-98=10936 zeroes. We draw with replacement a sample of 11037 items from the first population, and a sample of 11034 items from the second population. Each of these is called a *bootstrap sample*. From these we derive the bootstrap replicate of $\hat{\theta}$:

$$\hat{\theta}^* = \frac{\text{Proportion of ones in bootstrap sample } \#1}{\text{Proportion of ones in bootstrap sample } \#2}.$$
 (1.6)

We repeat this process a large number of times, say 1000 times, and obtain 1000 bootstrap replicates $\hat{\theta}^*$. This process is easy to implement on a computer, as we will see later. These 1000 replicates contain information that can be used to make inferences from our data. For example, the standard deviation turned out to be 0.17 in a batch of 1000 replicates that we generated. The value 0.17 is an estimate of the standard error of the ratio of rates $\hat{\theta}$. This indicates that the observed ratio $\hat{\theta} = 1.21$ is only a little more than one standard error larger than 1, and so the neutral value $\theta = 1$ cannot be ruled out. A rough 95% confidence interval like (1.5) can be derived by taking the 25th and 975th largest of the 1000 replicates, which in this case turned out to be (.93, 1.60).

In this simple example, the confidence interval derived from the bootstrap agrees very closely with the one derived from statistical theory. Bootstrap methods are intended to simplify the calculation of inferences like (1.2) and (1.5), producing them in an automatic way even in situations much more complicated than the aspirin study.

The terminology of statistical summaries and inferences, like regression, correlation, analysis of variance, discriminant analysis, standard error, significance level and confidence interval, has become the lingua franca of all disciplines that deal with noisy data. We will be examining what this language means and how it works in practice. The particular goal of bootstrap theory is a computerbased implementation of basic statistical concepts. In some ways it is easier to understand these concepts in computer-based contexts than through traditional mathematical exposition.

1.1 An overview of this book

This book describes the bootstrap and other methods for assessing statistical accuracy. The bootstrap does not work in isolation but rather is applied to a wide variety of statistical procedures. Part of the objective of this book is expose the reader to many exciting and useful statistical techniques through real-data examples. Some of the techniques described include nonparametric regression, density estimation, classification trees, and least median of squares regression.

Here is a chapter-by-chapter synopsis of the book. Chapter 2 introduces the bootstrap estimate of standard error for a simple mean. Chapters 3-5 contain some basic background material, and may be skimmed by readers eager to get to the details of the bootstrap in Chapter 6. Random samples, populations, and basic probability theory are reviewed in Chapter 3. Chapter 4 defines the empirical distribution function estimate of the population, which simply estimates the probability of each of n data items to be 1/n. Chapter 4 also shows that many familiar statistics can be viewed as "plug-in" estimates, that is, estimates obtained by plugging in the empirical distribution function for the unknown distribution of the population. Chapter 5 reviews standard error estimation for a mean, and shows how the usual textbook formula can be derived as a simple plug-in estimate.

The bootstrap is defined in **Chapter 6**, for estimating the standard error of a statistic from a single sample. The bootstrap standard error estimate is a plug-in estimate that rarely can be computed exactly; instead a simulation ("resampling") method is used for approximating it.

Chapter 7 describes the application of bootstrap standard errors in two complicated examples: a principal components analysis

and a curve fitting problem.

Up to this point, only one-sample data problems have been discussed. The application of the bootstrap to more complicated data structures is discussed in **Chapter 8**. A two-sample problem and a time-series analysis are described.

Regression analysis and the bootstrap are discussed and illustrated in **Chapter 9**. The bootstrap estimate of standard error is applied in a number of different ways and the results are discussed in two examples.

The use of the bootstrap for estimation of bias is the topic of **Chapter 10**, and the pros and cons of bias correction are discussed. **Chapter 11** describes the jackknife method in some detail. We see that the jackknife is a simple closed-form approximation to the bootstrap, in the context of standard error and bias estimation.

The use of the bootstrap for construction of confidence intervals is described in **Chapters 12**, **13** and **14**. There are a number of different approaches to this important topic and we devote quite a bit of space to them. In **Chapter 12** we discuss the bootstrap-tapproach, which generalizes the usual Student's t method for constructing confidence intervals. The percentile method (**Chapter 13**) uses instead the percentiles of the bootstrap distribution to define confidence limits. The BC_a (bias-corrected accelerated interval) makes important corrections to the percentile interval and is described in **Chapter 14**.

Chapter 15 covers permutation tests, a time-honored and useful set of tools for hypothesis testing. Their close relationship with the bootstrap is discussed; Chapter 16 shows how the bootstrap can be used in more general hypothesis testing problems.

Prediction error estimation arises in regression and classification problems, and we describe some approaches for it in **Chapter 17**. Cross-validation and bootstrap methods are described and illustrated. Extending this idea, **Chapter 18** shows how the bootstrap and cross-validation can be used to adapt estimators to a set of data.

Like any statistic, bootstrap estimates are random variables and so have inherent error associated with them. When using the bootstrap for making inferences, it is important to get an idea of the magnitude of this error. In **Chapter 19** we discuss the jackknifeafter-bootstrap method for estimating the standard error of a bootstrap quantity.

Chapters 20-25 contain more advanced material on selected

topics, and delve more deeply into some of the material introduced in the previous chapters. The relationship between the bootstrap and jackknife is studied via the "resampling picture" in **Chapter 20. Chapter 21** gives an overview of non-parametric and parametric inference, and relates the bootstrap to a number of other techniques for estimating standard errors. These include the delta method, Fisher information, infinitesimal jackknife, and the sandwich estimator.

Some advanced topics in bootstrap confidence intervals are discussed in Chapter 22, providing some of the underlying basis for the techniques introduced in Chapters 12–14. Chapter 23 describes methods for efficient computation of bootstrap estimates including control variates and importance sampling. In Chapter 24 the construction of approximate likelihoods is discussed. The bootstrap and other related methods are used to construct a "nonparametric" likelihood in situations where a parametric model is not specified.

Chapter 25 describes in detail a bioequivalence study in which the bootstrap is used to estimate power and sample size. In Chapter 26 we discuss some general issues concerning the bootstrap and its role in statistical inference.

Finally, the **Appendix** contains a description of a number of different computer programs for the methods discussed in this book.

1.2 Information for instructors

We envision that this book can provide the basis for (at least) two different one semester courses. An upper-year undergraduate or first-year graduate course could be taught from some or all of the first 19 chapters, possibly covering Chapter 25 as well (both authors have done this). In addition, a more advanced graduate course could be taught from a selection of Chapters 6–19, and a selection of Chapters 20–26. For an advanced course, supplementary material might be used, such as Peter Hall's book *The Bootstrap and Edgeworth Expansion* or journal papers on selected technical topics. The Bibliographic notes in the book contain many suggestions for background reading.

We have provided numerous exercises at the end of each chapter. Some of these involve computing, since it is important for the student to get hands-on experience for learning the material. The bootstrap is most effectively used in a high-level language for data analysis and graphics. Our language of choice (at present) is "S" (or "S-PLUS"), and a number of S programs appear in the Appendix. Most of these programs could be easily translated into other languages such as Gauss, Lisp-Stat, or Matlab. Details on the availability of S and S-PLUS are given in the Appendix.

1.3 Some of the notation used in the book

Lower case bold letters such as x refer to vectors, that is, $\mathbf{x} = (x_1, x_2, \dots, x_n)$. Matrices are denoted by upper case bold letters such as X, while a plain uppercase letter like X refers to a random variable. The transpose of a vector is written as \mathbf{x}^T . A superscript "*" indicates a bootstrap random variable: for example, \mathbf{x}^* indicates a bootstrap data set generated from a data set x. Parameters are denoted by Greek letters such as θ . A hat on a letter indicates an estimate, such as $\hat{\theta}$. The letters F and G refer to populations. In Chapter 21 the same symbols are used for the cumulative distribution function of a population. I_C is the indicator function equal to 1 if condition C is true and 0 otherwise. For example, $I_{\{x<2\}} = 1$ if x < 2 and 0 otherwise. The notation tr(A) refers to the trace of the matrix A, that is, the sum of the diagonal elements. The derivatives of a function g(x) are denoted by g'(x), g''(x) and so on.

The notation

$$F \rightarrow (x_1, x_2, \ldots x_n)$$

indicates an independent and identically distributed sample drawn from F. Equivalently, we also write $x_i \stackrel{\text{i.i.d.}}{\sim} F$ for i = 1, 2, ..., n.

Notation such as $\#\{x_i > 3\}$ means the number of x_i s greater than 3. log x refers to the natural logarithm of x.

CHAPTER 2

The accuracy of a sample mean

The bootstrap is a computer-based method for assigning measures of accuracy to statistical estimates. The basic idea behind the bootstrap is very simple, and goes back at least two centuries. After reviewing some background material, this book describes the bootstrap method, its implementation on the computer, and its application to some real data analysis problems. First though, this chapter focuses on the one example of a statistical estimator where we really don't need a computer to assess accuracy: the sample mean. In addition to previewing the bootstrap, this gives us a chance to review some fundamental ideas from elementary statistics. We begin with a simple example concerning means and their estimated accuracies.

Table 2.1 shows the results of a small experiment, in which 7 out of 16 mice were randomly selected to receive a new medical treatment, while the remaining 9 were assigned to the non-treatment (control) group. The treatment was intended to prolong survival after a test surgery. The table shows the survival time following surgery, in days, for all 16 mice.

Did the treatment prolong survival? A comparison of the means for the two groups offers preliminary grounds for optimism. Let x_1, x_2, \dots, x_7 indicate the lifetimes in the treatment group, so $x_1 =$ $94, x_2 = 197, \dots, x_7 = 23$, and likewise let y_1, y_2, \dots, y_9 indicate the control group lifetimes. The group means are

$$\bar{x} = \sum_{i=1}^{7} x_i/7 = 86.86$$
 and $\bar{y} = \sum_{i=1}^{9} y_i/9 = 56.22,$ (2.1)

so the difference $\bar{x} - \bar{y}$ equals 30.63, suggesting a considerable lifeprolonging effect for the treatment.

But how accurate are these estimates? After all, the means (2.1) are based on small samples, only 7 and 9 mice, respectively. In

Group	Data			(Sample Size)	Mean	Estimated Standard Error	
Treatment:	94	197	16				
	38	99	141				
	23			(7)	86.86	25.24	
Control:	52	104	146				
	10	51	30				
	40	27	46	(9)	56.22	14.14	
				Difference:	30.63	28.93	

Table 2.1. The mouse data. Sixteen mice were randomly assigned to a treatment group or a control group. Shown are their survival times, in days, following a test surgery. Did the treatment prolong survival?

order to answer this question, we need an estimate of the accuracy of the sample means \bar{x} and \bar{y} . For sample means, and essentially only for sample means, an accuracy formula is easy to obtain.

The estimated standard error of a mean \bar{x} based on n independent data points x_1, x_2, \dots, x_n , $\bar{x} = \sum_{i=1}^n x_i/n$, is given by the formula

$$\sqrt{\frac{s^2}{n}} \tag{2.2}$$

where $s^2 = \sum_{i=1}^{n} (x_i - \bar{x})^2 / (n-1)$. (This formula, and standard errors in general, are discussed more carefully in Chapter 5.) The standard error of any estimator is defined to be the square root of its variance, that is, the estimator's root mean square variability around its expectation. This is the most common measure of an estimator's accuracy. Roughly speaking, an estimator will be less than one standard error away from its expectation about 68% of the time, and less than two standard errors away about 95% of the time.

If the estimated standard errors in the mouse experiment were very small, say less than 1, then we would know that \bar{x} and \bar{y} were close to their expected values, and that the observed difference of 30.63 was probably a good estimate of the true survival-prolonging capability of the treatment. On the other hand, if formula (2.2) gave big estimated standard errors, say 50, then the difference estimate would be too inaccurate to depend on.

The actual situation is shown at the right of Table 2.1. The estimated standard errors, calculated from (2.2), are 25.24 for \bar{x} and 14.14 for \bar{y} . The standard error for the difference $\bar{x} - \bar{y}$ equals $28.93 = \sqrt{25.24^2 + 14.14^2}$ (since the variance of the difference of two independent quantities is the sum of their variances). We see that the observed difference 30.63 is only 30.63/28.93 = 1.05 estimated standard errors greater than zero. Readers familiar with hypothesis testing theory will recognize this as an *insignificant* result, one that could easily arise by chance even if the treatment really had no effect at all.

There are more precise ways to verify this disappointing result, (e.g. the permutation test of Chapter 15), but usually, as in this case, estimated standard errors are an excellent first step toward thinking critically about statistical estimates. Unfortunately standard errors have a major disadvantage: for most statistical estimators other than the mean there is no formula like (2.2) to provide estimated standard errors. In other words, it is hard to assess the accuracy of an estimate other than the mean.

Suppose for example, we want to compare the two groups in Table 2.1 by their medians rather than their means. The two medians are 94 for treatment and 46 for control, giving an estimated difference of 48, considerably more than the difference of the means. But how accurate are these medians? Answering such questions is where the bootstrap, and other computer-based techniques, come in. The remainder of this chapter gives a brief preview of the bootstrap estimate of standard error, a method which will be fully discussed in succeeding chapters.

Suppose we observe independent data points x_1, x_2, \dots, x_n , for convenience denoted by the vector $\mathbf{x} = (x_1, x_2, \dots, x_n)$, from which we compute a statistic of interest $s(\mathbf{x})$. For example the data might be the n = 9 control group observations in Table 2.1, and $s(\mathbf{x})$ might be the sample mean.

The bootstrap estimate of standard error, invented by Efron in 1979, looks completely different than (2.2), but in fact it is closely related, as we shall see. A *bootstrap sample* $\mathbf{x}^* = (x_1^*, x_2^*, \dots, x_n^*)$ is obtained by randomly sampling n times, with replacement, from the original data points x_1, x_2, \dots, x_n . For instance, with n = 7 we might obtain $\mathbf{x}^* = (x_5, x_7, x_5, x_4, x_7, x_3, x_1)$.



Figure 2.1. Schematic of the bootstrap process for estimating the standard error of a statistic $s(\mathbf{x})$. B bootstrap sample, are generated from the original data set. Each bootstrap sample has n elements, generated by sampling with replacement n times from the original data set. Bootstrap replicates $s(\mathbf{x}^{*1}), s(\mathbf{x}^{*2}), \ldots s(\mathbf{x}^{*B})$ are obtained by calculating the value of the statistic $s(\mathbf{x})$ on each bootstrap sample. Finally, the standard deviation of the values $s(\mathbf{x}^{*1}), s(\mathbf{x}^{*2}), \ldots s(\mathbf{x}^{*B})$ is our estimate of the standard error of $s(\mathbf{x})$.

Figure 2.1 is a schematic of the bootstrap process. The bootstrap algorithm begins by generating a large number of independent bootstrap samples $\mathbf{x}^{*1}, \mathbf{x}^{*2}, \dots, \mathbf{x}^{*B}$, each of size *n*. Typical values for *B*, the number of bootstrap samples, range from 50 to 200 for standard error estimation. Corresponding to each bootstrap sample is a *bootstrap replication* of *s*, namely $s(\mathbf{x}^{*b})$, the value of the statistic *s* evaluated for \mathbf{x}^{*b} . If $s(\mathbf{x})$ is the sample median, for instance, then $s(\mathbf{x}^*)$ is the median of the bootstrap sample. The bootstrap estimate of standard error is the standard deviation of the bootstrap replications,

$$\widehat{se}_{boot} = \left\{ \sum_{b=1}^{B} [s(\mathbf{x}^{*b}) - s(\cdot)]^2 / (B-1) \right\}^{\frac{1}{2}},$$
(2.3)

where $s(\cdot) = \sum_{b=1}^{B} s(\mathbf{x}^{*b})/B$. Suppose $s(\mathbf{x})$ is the mean \bar{x} . In this

Table 2.2. Bootstrap estimates of standard error for the mean and median; treatment group, mouse data, Table 2.1. The median is less accurate (has larger standard error) than the mean for this data set.

B:	50	100	250	500	1000	∞
mean:	19.72	23.63	22.32	23.79	23.02	23.36
median:	32.21	36.35	34.46	36.72	36.48	37.83

case, standard probability theory tells us (Problem 2.5) that as B gets very large, formula (2.3) approaches

$$\{\sum_{i=1}^{n} (x_i - \bar{x})^2 / n^2\}^{\frac{1}{2}}.$$
(2.4)

This is almost the same as formula (2.2). We could make it exactly the same by multiplying definition (2.3) by the factor $[n/(n-1)]^{\frac{1}{2}}$, but there is no real advantage in doing so.

Table 2.2 shows bootstrap estimated standard errors for the mean and the median, for the treatment group mouse data of Table 2.1. The estimated standard errors settle down to limiting values as the number of bootstrap samples B increases. The limiting value 23.36 for the mean is obtained from (2.4). The formula for the limiting value 37.83 for the standard error of the median is quite complicated: see Problem 2.4 for a derivation.

We are now in a position to assess the precision of the difference in medians between the two groups. The bootstrap procedure described above was applied to the control group, producing a standard error estimate of 11.54 based on B = 100 replications ($B = \infty$ gave 9.73). Therefore, using B = 100, the observed difference of 48 has an estimated standard error of $\sqrt{36.35^2 + 11.54^2} = 38.14$, and hence is 48/38.14 = 1.26 standard errors greater than zero. This is larger than the observed difference in means, but is still insignificant.

For most statistics we don't have a formula for the limiting value of the standard error, but in fact no formula is needed. Instead we use the numerical output of the bootstrap program, for some convenient value of B. We will see in Chapters 6 and 19, that Bin the range 50 to 200 usually makes \hat{se}_{boot} a good standard error

PROBLEMS

estimator, even for estimators like the median. It is easy to write a bootstrap program that works for any computable statistic $s(\mathbf{x})$, as shown in Chapters 6 and the Appendix. With these programs in place, the data analyst is free to use any estimator, no matter how complicated, with the assurance that he or she will also have a reasonable idea of the estimator's accuracy. The price, a factor of perhaps 100 in increased computation, has become affordable as computers have grown faster and cheaper.

Standard errors are the simplest measures of statistical accuracy. Later chapters show how bootstrap methods can assess more complicated accuracy measures, like biases, prediction errors, and confidence intervals. Bootstrap confidence intervals add another factor of 10 to the computational burden. The payoff for all this computation is an increase in the statistical problems that can be analyzed, a reduction in the assumptions of the analysis, and the elimination of the routine but tedious theoretical calculations usually associated with accuracy assessment.

2.1 Problems

- 2.1[†] Suppose that the mouse survival times were expressed in weeks instead of days, so that the entries in Table 2.1 were all divided by 7.
 - (a) What effect would this have on \bar{x} and on its estimated standard error (2.2)? Why does this make sense?
 - (b) What effect would this have on the ratio of the difference $\bar{x} \bar{y}$ to its estimated standard error?
- 2.2 Imagine the treatment group in Table 2.1 consisted of R repetitions of the data actually shown, where R is a positive integer. That is, the treatment data consisted of R 94's, R 197's, etc. What effect would this have on the estimated standard error (2.2)?
- 2.3 It is usually true that the error of a statistical estimator decreases at a rate of about 1 over the square root of the sample size. Does this agree with the result of Problem 2.2?
- 2.4 Let $x_{(1)} < x_{(2)} < x_{(3)} < x_{(4)} < x_{(5)} < x_{(6)} < x_{(7)}$ be an ordered sample of size n = 7. Let \mathbf{x}^* be a bootstrap sample, and $s(\mathbf{x}^*)$ be the corresponding bootstrap replication of the median. Show that

- (a) $s(\mathbf{x}^*)$ equals one of the original data values $x_{(i)}$, $i = 1, 2, \dots, 7$.
- (b) [†] $s(\mathbf{x}^*)$ equals $x_{(i)}$ with probability

$$p(i) = \sum_{j=0}^{3} \{ \operatorname{Bi}(j; n, \frac{i-1}{n}) - \operatorname{Bi}(j; n, \frac{i}{n}) \},$$
(2.5)

where Bi(j; n, p) is the binomial probability $\binom{n}{j}p^{j}(1-p)^{n-j}$. [The numerical values of p(i) are .0102, .0981, .2386, .3062, .2386, .0981, .0102. These values were used to compute \widehat{se}_{boot} { median} = 37.83, for $B = \infty$, Table 2.2.]

- 2.5 Apply the weak law of large numbers to show that expression (2.3) approaches expression (2.4) as n goes to infinity.
- † Indicates a difficult or more advanced problem.

CHAPTER 3

Random samples and probabilities

3.1 Introduction

Statistics is the theory of accumulating information, especially information that arrives a little bit at a time. A typical statistical situation was illustrated by the mouse data of Table 2.1. No one mouse provides much information, since the individual results are so variable, but seven, or nine mice considered together begin to be quite informative. Statistical theory concerns the best ways of extracting this information. Probability theory provides the mathematical framework for statistical inference. This chapter reviews the simplest probabilistic model used to model random data: the case where the observations are a random sample from a single unknown population, whose properties we are trying to learn from the observed data.

3.2 Random samples

It is easiest to visualize random samples in terms of a finite population or "universe" \mathcal{U} of individual units U_1, U_2, \dots, U_N , any one of which is equally likely to be selected in a single random draw. The population of units might be all the registered voters in an area undergoing a political survey, all the men that might conceivably be selected for a medical experiment, all the high schools in the United States, etc. The individual units have properties we would like to learn, like a political opinion, a medical survival time, or a graduation rate. It is too difficult and expensive to examine every unit in \mathcal{U} , so we select for observation a random sample of manageable size.

A random sample of size n is defined to be a collection of n

units u_1, u_2, \dots, u_n selected at random from \mathcal{U} . In principle the sampling process goes as follows: a random number device independently selects integers j_1, j_2, \dots, j_n , each of which equals any value between 1 and N with probability 1/N. These integers determine which members of \mathcal{U} are selected to be in the random sample, $u_1 = U_{j_1}, u_2 = U_{j_2}, \dots, u_n = U_{j_n}$. In practice the selection process is seldom this neat, and the population \mathcal{U} may be poorly defined, but the conceptual framework of random sampling is still useful for understanding statistical inference. (The methodology of good experimental design, for example the random assignment of selected units to Treatment or Control groups as was done in the mouse experiment, helps make random sampling theory more applicable to real situations like that of Table 2.1.)

Our definition of random sampling allows a single unit U_i to appear more than once in the sample. We could avoid this by insisting that the integers j_1, j_2, \dots, j_n be distinct, called "sampling without replacement." It is a little simpler to allow repetitions, that is to "sample with replacement", as in the previous paragraph. If the size n of the random sample is much smaller than the population size N, as is usually the case, the probability of sample repetitions will be small anyway. See Problem 3.1. Random sampling always means sampling with replacement in what follows, unless otherwise stated.

Having selected a random sample u_1, u_2, \dots, u_n , we obtain one or more measurements of interest for each unit. Let x_i indicate the measurements for unit u_i . The observed data are the collection of measurements x_1, x_2, \dots, x_n . Sometimes we will denote the observed data (x_1, x_2, \dots, x_n) by the single symbol **x**.

We can imagine making the measurements of interest on every member U_1, U_2, \dots, U_N of \mathcal{U} , obtaining values X_1, X_2, \dots, X_N . This would be called a census of U.

The symbol \mathcal{X} will denote the census of measurements (X_1, X_2, \cdots, X_N) . We will also refer to \mathcal{X} as the population of measurements, or simply the population, and call **x** a random sample of size n from \mathcal{X} . In fact, we usually can't afford to conduct a census, which is why we have taken a random sample. The goal of statistical inference is to say what we have learned about the population \mathcal{X} from the observed data **x**. In particular, we will use the bootstrap to say how accurately a statistic calculated from x_1, x_2, \cdots, x_n (for instance the sample median) estimates the corresponding quantity for the whole population.

Table 3.1. The law school data. A random sample of size n = 15 was taken from the collection of N = 82 American law schools participating in a large study of admission practices. Two measurements were made on the entering classes of each school in 1973: LSAT, the average score for the class on a national law test, and GPA, the average undergraduate grade-point average for the class.

School	LSAT	GPA	School	LSAT	GPA
1	576	3.39	9	651	3.36
2	635	3.30	10	605	3.13
3	558	2.81	11	653	3.12
4	578	3.03	12	575	2.74
5	666	3.44	13	545	2.76
6	580	3.07	14	572	2.88
7	555	3.00	15	594	2.96
8	661	3.43			

Table 3.1 shows a random sample of size n = 15 drawn from a population of N = 82 American law schools. What is actually shown are two measurements made on the entering classes of 1973 for each school in the sample: LSAT, the average score of the class on a national law test, and GPA, the average undergraduate grade point average achieved by the members of the class. In this case the measurement x_i on u_i , the *i*th member of the sample, is the pair

$$x_i = (\text{LSAT}_i, \text{GPA}_i)$$
 $i = 1, 2, \cdots, 15.$

The observed data x_1, x_2, \dots, x_n is the collection of 15 pairs of numbers shown in Table 3.1.

This example is an artificial one because the census of data X_1, X_2, \dots, X_{82} was actually made. In other words, LSAT and GPA are available for the entire population of N = 82 schools. Figure 3.1 shows the census data and the sample data. Table 3.2 gives the entire population of N measurements.

In a real statistical problem, like that of Table 3.1, we would see only the sample data, from which we would be trying to infer the properties of the population. For example, consider the 15 LSAT scores in the observed sample. These have mean 600.27 with estimated standard error 10.79, based on the data in Table 3.1 and formula (2.2). There is about a 68% chance that the true LSAT



Figure 3.1. The left panel is a scatterplot of the (LSAT, GPA) data for all N = 82 law schools; circles indicate the n = 15 data points comprising the "observed sample" of Table 3.1. The right panel shows only the observed sample. In problems of statistical inference, we are trying to infer the situation on the left from the picture on the right.

mean, the mean for the entire population from which the observed data was sampled, lies in the interval 600.27 ± 10.79 .

We can check this result, since we are dealing with an artificial example for which the complete population data are known. The mean of all 82 LSAT values is 597.55, lying nicely within the predicted interval 600.27 ± 10.79 .

3.3 Probability theory

Statistical inference concerns learning from experience: we observe a random sample $\mathbf{x} = (x_1, x_2, \dots, x_n)$ and wish to infer properties of the complete population $\mathcal{X} = (X_1, X_2, \dots, X_N)$ that yielded the sample. Probability theory goes in the opposite direction: from the composition of a population \mathcal{X} we deduce the properties of a random sample \mathbf{x} , and of statistics calculated from \mathbf{x} . Statistical inference as a mathematical science has been developed almost exclusively in terms of probability theory. Here we will review briefly

school	LSAT	GPA	school	LSAT	GPA	school	LSAT	GPA
1	622	3.23	28	632	3.29	56	641	3.28
2	542	2.83	29	587	3.16	57	512	3.01
3	579	3.24	30	581	3.17	58	631	3.21
4+	653	3.12	31+	605	3.13	59	597	3.32
5	606	3.09	32	704	3.36	60	621	3.24
6+	576	3.39	33	477	2.57	61	617	3.03
7	620	3.10	34	591	3.02	62	637	3.33
8	615	3.40	35+	578	3.03	62	572	3.08
9	553	2.97	36+	572	2.88	64	610	3.13
10	607	2.91	37	615	3.37	65	562	3.01
11	558	3.11	38	606	3.20	66	635	3.30
12	596	3.24	39	603	3.23	67	614	3.15
13+	635	3.30	40	535	2.98	68	546	2.82
14	581	3.22	41	595	3.11	69	598	3.20
15 +	661	3.43	42	575	2.92	70+	666	3.44
16	547	2.91	43	573	2.85	71	570	3.01
17	599	3.23	44	644	3.38	72	570	2.92
18	646	3.47	45+	545	2.76	73	605	3.45
19	622	3.15	46	645	3.27	74	565	3.15
20	611	3.33	47+	651	3.36	75	686	3.50
21	546	2.99	48	562	3.19	76	608	3.16
22	614	3.19	49	609	3.17	77	595	3.19
23	628	3.03	50+	555	3.00	78	590	3.15
24	575	3.01	51	586	3.11	79+	558	2.81
25	662	3.39	52+	580	3.07	80	611	3.16
26	627	3.41	53+	594	2.96	81	564	3.02
27	608	3.04	54	594	3.05	82+	575	2.74
			55	560	2.93			

Table 3.2. The population of measurements (LSAT, GPA), for the universe of 82 law schools. The data in Table 3.1 was sampled from this population. The +'s indicate the sampled schools.

some fundamental concepts of probability, including probability distributions, expectations, and independence.

As a first example, let x represent the outcome of rolling a fair die so x is equally likely to be 1, 2, 3, 4, 5, or 6. We write this in probability notation as

$$Prob\{x = k\} = 1/6 \quad \text{for} \quad k = 1, 2, 3, 4, 5, 6. \tag{3.1}$$

A random quantity like x is often called a *random variable*.

Probabilities are idealized or theoretical proportions. We can imagine a universe $\mathcal{U} = \{U_1, U_2, \dots, U_N\}$ of possible rolls of the die, where U_j completely describes the physical act of the *j*th roll, with corresponding results $\mathcal{X} = (X_1, X_2, \dots, X_N)$. Here N might be very large, or even infinite. The statement $\operatorname{Prob}\{x = 5\} = 1/6$ means that a randomly selected member of \mathcal{X} has a 1/6 chance of equaling 5, or more simply that 1/6 of the members of \mathcal{X} equal 5. Notice that probabilities, like proportions, can never be less than 0 or greater than 1.

For convenient notation define the frequencies f_k ,

$$f_k = \operatorname{Prob}\{x = k\},\tag{3.2}$$

so the fair die has $f_k = 1/6$ for $k = 1, 2, \dots, 6$. The probability distribution of a random variable x, which we will denote by F, is any complete description of the probabilistic behavior of x. F is also called the probability distribution of the population \mathcal{X} . Here we can take F to be the vector of frequencies

$$F = (f_1, f_2, \cdots, f_6) = (1/6, 1/6, \cdots, 1/6).$$
(3.3)

An unfair die would be one for which F did not equal $(1/6, 1/6, \dots, 1/6)$.

<u>Note</u>: In many books, the symbol F is used for the cumulative probability distribution function $F(x_0) = \operatorname{Prob}\{x \leq x_0\}$ for $-\infty < x_0 < \infty$. This is an equally valid description of the probabilistic behavior of x, but it is only convenient for the case where x is a real number. We will also be interested in cases where x is a vector, as in Table 3.1, or an even more general object. This is the reason for defining F as any description of x's probabilities, rather than the specific description in terms of the cumulative probabilities. When no confusion can arise, in later chapters we use symbols like F and G to represent cumulative distribution functions.

Some probability distributions arise so frequently that they have received special names. A random variable x is said to have the *binomial distribution* with size n and probability of success p, denoted

$$x \sim \operatorname{Bi}(n, p),$$
 (3.4)

if its frequencies are

$$f_k = \binom{n}{k} p^k (1-p)^{n-k} \quad \text{for} \quad k = 0, 1, 2, \cdots, n.$$
 (3.5)

Here n is a positive integer, p is a number between 0 and 1, and $\binom{n}{k}$ is the binomial coefficient n!/[k!(n-k)!]. Figure 3.2 shows the

distribution $F = (f_0, f_1, \dots, f_n)$ for $x \sim \text{Bi}(n, p)$, with n = 25 and p = .25, .50, and .90. We also write F = Bi(n, p) to indicate situation (3.4).

Let A be a set of integers. Then the probability that x takes a value in A, or more simply the probability of A, is

$$\operatorname{Prob}\{x \in A\} = \operatorname{Prob}\{A\} = \sum_{k \in A} f_k.$$
(3.6)

For example if $A = \{1, 3, 5, \dots, 25\}$ and $x \sim \text{Bi}(25, p)$, then $\text{Prob}\{A\}$ is the probability that a binomial random variable of size 25 and probability of success p equals an odd integer. Notice that since f_k is the theoretical proportion of times x equals k, the sum $\sum_{k \in A} f_k = \text{Prob}\{A\}$ is the theoretical proportion of times x takes its value in A.

The sample space of x, denoted S_x , is the collection of possible values x can have. For a fair die, $S_x = \{1, 2, \dots, 6\}$, while $S_x = \{0, 1, 2, \dots, n\}$ for a Bi(n, p) distribution. By definition, x occurs in S_x every time, that is, with theoretical proportion 1, so

$$\operatorname{Prob}\{\mathcal{S}_x\} = \sum_{k \in \mathcal{S}_x} f_k = 1.$$
(3.7)

For any probability distribution on the integers the frequencies f_j are nonnegative numbers summing to 1.

In our examples so far, the sample space S_x has been a subset of the integers. One of the convenient things about probability distributions is that they can be defined on quite general spaces. Consider the law school data of Figure 3.1. We might take S_x to be the positive quadrant of the plane,

$$S_x = \mathcal{R}^{2+} = \{(y, z), y > 0, z > 0\}.$$
(3.8)

(This includes values like $x = (10^6, 10^9)$, but it doesn't hurt to let S_x be too big.) For a subset A of S_x , we would still write Prob $\{A\}$ to indicate the probability that x occurs in A.

For example, we could take

$$A = \{(y, z) : 0 < y < 600, 0 < z < 3.0\}.$$
(3.9)

A law school $x \in A$ if its 1973 entering class had LSAT less than 600 and GPA less than 3.0. In this case we happen to know the complete population \mathcal{X} ; it is the 82 points indicated on the left panel of Figure 3.1 and in Table 3.2. Of these, 16 are in A, so

$$Prob\{A\} = 16/82 = .195. \tag{3.10}$$



Figure 3.2. The frequencies f_0, f_1, \dots, f_n for the binomial distributions Bi(n, p), n = 25 and p = .25, .50, and .90. The points have been connected by lines to enhance visibility.

Here the idealized proportion $\operatorname{Prob}\{A\}$ is an actual proportion. Only in cases where we have a complete census of the population is it possible to directly evaluate probabilities as proportions.

The probability distribution F of x is still defined to be any complete description of x's probabilities. In the law school example, F can be described as follows: for any subset A of $S_x = \mathcal{R}^{2+}$,

$$Prob\{x \in A\} = \#\{X_j \in A\}/82, \tag{3.11}$$

where $\#\{X_j \in A\}$ is the number of the 82 points in the left panel of Figure 3.1 that lie in A. Another way to say the same thing is that F is a discrete distribution putting probability (or frequency) 1/82 on each of the indicated 82 points.

Probabilities can be defined continuously, rather than discretely as in (3.6) or (3.11). The most famous example is the *normal* (or *Gaussian*, or *bell-shaped*) distribution. A real-valued random variable x is defined to have the normal distribution with mean μ and variance σ^2 , written

$$x \sim N(\mu, \sigma^2)$$
 or $F = N(\mu, \sigma^2)$, (3.12)

if

$$\operatorname{Prob}\{x \in A\} = \int_{A} \frac{1}{\sqrt{2\pi\sigma^{2}}} e^{-\frac{1}{2}(\frac{x-\mu}{\sigma})^{2}} dx$$
(3.13)

for any subset A of the real line \mathcal{R}^1 . The integral in (3.13) is over the values of $x \in A$.

There are higher dimensional versions of the normal distribution, which involve taking integrals similar to (3.13) over multidimensional sets A. We won't need continuous distributions for development of the bootstrap (though they will appear later in some of the applications) and will avoid mathematical derivations based on calculus. As we shall see, one of the main incentives for the development of the bootstrap is the desire to substitute computer power for theoretical calculations involving special distributions.

The expectation of a real-valued random variable x, written E(x), is its average value, where the average is taken over the possible outcomes of x weighted according to its probability distribution F. Thus

$$E(x) = \sum_{x=0}^{n} x \binom{n}{x} p^{x} (1-p)^{x} \text{ for } x \sim Bi(n,p), \qquad (3.14)$$

and

$$E(x) = \int_{-\infty}^{\infty} x \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{1}{2}(\frac{x-\mu}{\sigma})^2} dx \text{ for } x \sim N(\mu, \sigma^2). \quad (3.15)$$

It is not difficult to show that E(x) = np for $x \sim Bi(n, p)$, and $E(x) = \mu$ for $x \sim N(\mu, \sigma^2)$. (See Problems 3.6 and 3.7.)

We sometimes write the expectation as $E_F(x)$, to indicate that the average is taken with respect to the distribution F.

Suppose r = g(x) is some function of the random variable x. Then E(r), the expectation of r, is the theoretical average of g(x) weighted according to the probability distribution of x. For example if $x \sim N(\mu, \sigma^2)$ and $r = x^3$, then

$$\mathbf{E}(r) = \int_{-\infty}^{\infty} x^3 \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{1}{2}(\frac{x-\mu}{\sigma})^2} dx.$$
(3.16)

Probabilities are a special case of expectations. Let A be a subset

of S_x , and take $r = I_{\{x \in A\}}$ where $I_{\{x \in A\}}$ is the indicator function

$$I_{\{x \in A\}} = \begin{cases} 1 & \text{if } x \in A \\ 0 & \text{if } x \notin A \end{cases}.$$
(3.17)

Then E(r) equals $\operatorname{Prob}\{x \in A\}$, or equivalently

$$E(I_{\{x \in A\}}) = Prob\{x \in A\}.$$
 (3.18)

For example if $x \sim N(\mu, \sigma^2)$, then

$$E(r) = \int_{-\infty}^{\infty} I_{\{x \in A\}} \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{1}{2}(\frac{x-\mu}{\sigma})^2} dx$$

= $\int_{A} \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{1}{2}(\frac{x-\mu}{\sigma})^2} dx,$ (3.19)

which is $\operatorname{Prob}\{x \in A\}$ according to (3.13).

The notion of an expectation as a theoretical average is very general, and includes cases where the random variable x is not real-valued. In the law school situation, for instance, we might be interested in the expectation of the ratio of LSAT and GPA. Writing x = (y, z) as in (3.8), then r = y/z, and the expectation of r is

$$E(LSAT/GPA) = \frac{1}{82} \sum_{j=1}^{82} (y_j/z_j)$$
(3.20)

where $x_j = (y_j, z_j)$ is the *j*th point in Table 3.2. Numerical evaluation of (3.20) gives E(LSAT/GPA) = 190.8.

Let $\mu_x = E_F(x)$, for x a real-valued random variable with distribution F. The variance of x, indicated by σ_x^2 or just σ^2 , is defined to be the expected value of $y = (x - \mu)^2$. In other words, σ^2 is the theoretical average squared distance of a random variable x from its expectation μ_x ,

$$\sigma_x^2 = \mathcal{E}_F (x - \mu_x)^2. \tag{3.21}$$

The variance of $x \sim N(\mu, \sigma^2)$ equals σ^2 ; the variance of $x \sim \text{Bi}(n, p)$ equals np(1-p), see Problem 3.9. The standard deviation of a random variable is defined to be the square root of its variance.

Two random variables y and z are said to be *independent* if

$$\mathbf{E}[g(y)h(z)] = \mathbf{E}[g(y)]\mathbf{E}[h(z)]$$
(3.22)

for all functions g(y) and h(z). Independence is well named: (3.22) implies that the random outcome of y doesn't affect the random outcome of z, and vice-versa.

To see this, let B and C be subsets of S_y and S_z respectively, the sample spaces of y and z, and take g and h to be the indicator functions $g(y) = I_{\{y \in B\}}$ and $h(z) = I_{\{z \in C\}}$. Notice that

$$I_{\{y\in B\}}I_{\{z\in C\}} = \begin{cases} 1 & \text{if } y\in B \text{ and } z\in C\\ 0 & \text{otherwise.} \end{cases}$$
(3.23)

So $I_{\{y \in B\}}I_{\{z \in C\}}$ is the indicator function of the intersection $\{y \in B\} \cap \{z \in C\}$. Then by (3.18) and the independence definition (3.22),

$$\begin{aligned} \operatorname{Prob}\{(y,z) \in B \cap C\} &= \operatorname{E}(I_{\{y \in B\}}I_{\{z \in C\}}) = \operatorname{E}(I_{\{y \in B\}})\operatorname{E}(I_{\{z \in C\}}) \\ &= \operatorname{Prob}\{y \in B\}\operatorname{Prob}\{z \in C\}. \end{aligned}$$

$$(3.24)$$

Looking at Figure 3.1, we can see that (3.24) does *not* hold for the law school example, see Problem 3.10, so LSAT and GPA are not independent.

Whether or not y and z are independent, expectations follow the simple addition rule

$$E[g(y) + h(z)] = E[g(y)] + E[h(z)].$$
(3.25)

In general,

$$E[\sum_{i=1}^{n} g_i(x_i)] = \sum_{i=1}^{n} E[g_i(x_i)]$$
(3.26)

for any functions g_i of any *n* random variables x_1, x_2, \dots, x_n .

Random sampling with replacement guarantees independence: if $\mathbf{x} = (x_1, x_2, \dots, x_n)$ is a random sample of size n from a population \mathcal{X} , then all n observations x_i are identically distributed and mutually independent of each other. In other words, all of the x_i have the same probability distribution F, and

for any functions g_1, g_2, \dots, g_n . (This is almost a definition of what random sampling means.) We will write

$$F \to (x_1, x_2, \cdots, x_n) \tag{3.28}$$

to indicate that $\mathbf{x} = (x_1, x_2, \dots, x_n)$ is a random sample of size n from a population with probability distribution F. This is sometimes written as

$$x_i \stackrel{\text{i.i.d.}}{\sim} F \qquad i = 1, 2, \cdots, n, \tag{3.29}$$

where i.i.d. stands for independent and identically distributed.

3.4 Problems

3.1 A random sample of size n is taken with replacement from a population of size N. Show that the probability of having no repetitions in the sample is given by the product

$$\prod_{j=0}^{n-1} (1-\frac{j}{N}).$$

- 3.2 Why might you suspect that the sample of 15 law schools in Table (3.1) was obtained by sampling without replacement, rather than with replacement?
- 3.3 The mean GPA for all 82 law schools is 3.13. How does this compare with the mean GPA for the observed sample of 15 law schools in Table 3.1? Is this difference compatible with the estimated standard error (2.2)?
- 3.4 Denote the mean and standard deviation of a set of numbers X_1, X_2, \dots, X_N by \overline{X} and S respectively, where

$$\overline{X} = \sum_{j=1}^{N} X_j / N$$
 $S = \{\sum_{j=1}^{N} (X_j - \overline{X})^2 / N\}^{1/2}$

(a) A sample x_1, x_2, \dots, x_n is selected from X_1, X_2, \dots, X_N by random sampling with replacement. Denote the standard deviation of the sample average $\bar{x} = \sum_{i=1}^n x_i/n$, usually called *the standard error* of \bar{x} , by $\operatorname{se}(\bar{x})$. Use a basic result of probability theory to show that

$$\operatorname{se}(\bar{x}) = \frac{S}{\sqrt{n}}.$$

(b) [†] Suppose instead that x_1, x_2, \dots, x_n is selected by random sampling *without* replacement (so we must have

 $n \leq N$), show that

$$\operatorname{se}(\bar{x}) = \frac{S}{\sqrt{n}} \left[\frac{N-n}{N-1} \right]^{\frac{1}{2}}$$

- (c) We see that sampling without replacement gives a smaller standard error for \bar{x} . Proportionally how much smaller will it be in the case of the law school data?
- 3.5 Given a random sample x_1, x_2, \dots, x_n , the *empirical proba*bility of a set A is defined to be the proportion of the sample in A, written

$$Prob{A} = \#{x_i \in A}/n.$$
 (3.30)

- (a) Find $\widehat{\text{Prob}}\{A\}$ for the data in Table 3.1, with A as given in (3.9).
- (b) The standard error of an empirical probability is $[\operatorname{Prob}\{A\} \cdot (1 \operatorname{Prob}\{A\})/n]^{1/2}$. How many standard errors is $\operatorname{Prob}\{A\}$ from $\operatorname{Prob}\{A\}$, given in (3.10)?
- 3.6 A very simple probability distribution F puts probability on only two outcomes, 0 or 1, with frequencies

$$f_0 = 1 - p, \quad f_1 = p.$$
 (3.31)

This is called the *Bernoulli* distribution. Here p is a number between 0 and 1. If x_1, \dots, x_n is a random sample from F, then elementary probability theory tells us that the sum

$$s = x_1 + x_2 + \dots + x_n$$
 (3.32)

has the binomial distribution (3.5),

$$s \sim \operatorname{Bi}(n, p).$$
 (3.33)

(a) Show that the empirical probability (3.30) satisfies

$$n \cdot \operatorname{Prob}\{A\} \sim \operatorname{Bi}(n, \operatorname{Prob}\{A\}).$$
 (3.34)

Expression (3.34) can also be written as $\widehat{\text{Prob}}\{A\} \sim \operatorname{Bi}(n, \operatorname{Prob}\{A\})/n.$

(b) Prove that if
$$x \sim Bi(n, p)$$
, then $E(x) = np$.

3.7 Without using calculus, give a symmetry argument to show that $E(x) = \mu$ for $x \sim N(\mu, \sigma^2)$.