FUNDAMENTALS OF

# Biostatistics



Bernard Rosne



3 1121 Dondosed



## FUNDAMENTALS OF BIOSTATISTICS

# BERNARD ROSNER HARVARD UNIVERSITY



#### **Duxbury Press**

An Imprint of Wadsworth Publishing Company An International Thomson Publishing Company

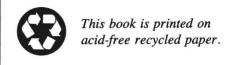
### This book is dedicated to my wife Cynthia and my children Sarah, David, and Laura.

Statistics Editor: Alexander Kugushev Editorial Assistant: Jennifer Burger Production Editor: The Book Company

Print Buyer: Karen Hunt Copy Editor: Linda Purrington Technical Illustrator: TUCKERdesign

Cover: William Reuter

Compositor: Jonathan Peck Typographers Printer: Banta Co., Harrisonburg, VA



#### COPYRIGHT © 1995

By Wadsworth Publishing Company A Division of International Thomson Publishing Inc.



Printed in the United States of America

For more information, contact:

Wadsworth Publishing
Company
1120 Birchmount Road
10 Davis Drive
Belmont, California 94002
International Thomson
Publishing
Nelson Canada
Canada M1K 5G4
International Thomson
Publishing
Publishing GmbH

Berkshire House 168-173

High Holborn

London, WC1V7AA

Königwinterer Strasse 418

53227 Bonn

Germany

England International Thomson
Thomas Nelson Australia Publishing Asia
102 Dodds Street 221 Henderson Road

South Melbourne 3205 #05-10 Victoria, Australia Singapore 0315 International Thomson Publishing-Japan Hirakawacho-cho Kyowa Building, 3F 2-2-1 Hirakawacho-cho Chiyoda-ku, 102 Tokyo Japan

All rights reserved. No part of this work covered by the copyright hereon may be reproduced or used in any form or by any means—graphic, electronic, or mechanical, including photocopying, recording, taping, or information storage and retrieval systems—without the written permission of the publisher.

4 5 6 7 8 9 10-01 00 99 98 97

#### Library of Congress Cataloging-in-Publication Data

Rosner, Bernard (Bernard A.)

Fundamentals of biostatistics / Bernard Rosner, -4th ed.

p. cm.

Includes bibliographical references and index.

ISBN 0-531-20840-8

1. Biometry. 2. Medical statistics. I. Title

OH323.5.R674 1994

574'.01'5185-dc20

### **PREFACE**

I have written this introductory-level biostatistics text for upper-level undergraduate or graduate students interested in medicine or other health-related areas. This book requires no previous background in statistics, and its mathematical level assumes only a knowledge of algebra.

Fundamentals of Biostatistics evolved from a set of notes that I used in a course in biostatistics taught to Harvard University undergraduates and Harvard Medical School students over the past fifteen years. I wrote this book to help motivate students to master the statistical methods that are most often used in the medical literature. From the student's viewpoint, it is important that the example material used to develop these methods is representative of what actually exists in the literature. Therefore, most examples and exercises used in this book are either based on actual articles from the medical literature or on actual medical research problems I have encountered during my consulting experience at the Harvard Medical School.

#### The Approach

Most other introductory statistics texts either use a completely nonmathematical, cookbook approach or develop the material in a rigorous, sophisticated mathematical framework. In this book I have attempted to follow an intermediate course, minimizing the amount of mathematical formulation and yet giving complete explanations of all the important concepts. *Every* new concept is developed systematically through completely worked out examples from current medical research problems. In addition, computer output is introduced where appropriate to illustrate these concepts.

The material in this book is suitable for either a one- or two-semester course in biostatistics. The material in Chapters 1 through 8 and Chapter 10 is suitable for a one-semester course. The instructor may select appropriate material from the other chapters as time permits.

#### Changes in the Fourth Edition

There are a total of 21 new sections and 9 additional sections with substantial revisions in the Fourth Edition. The new features include:

- An expanded set of computer exercises based on real data sets has been developed. The data sets are on a diskette that is provided with the book.
- A case study on lead exposure in children used in several chapters throughout the book.
- An extended discussion of randomized clinical trials including
  - (a) design features (Section 6.4.1)
  - (b) sample size issues (Section 10.7.3)
- One-Sample Inference for the Poisson distribution (Section 6.8 and 7.11)
- Sample size estimation based on confidence interval width (Section 7.7.3)
- Outlier detection techniques (Section 8.9)
- The one-way ANOVA random effects model (Section 9.6)

- The cross-over design (Section 9.7)
- A discussion of the most popular study designs in biomedical research (Section 10.3)
- Measures of effect for categorical data (Section 10.4)
- The hypergeometric distribution (Section 10.5.1)
- Issues in epidemiologic research include confounding, standardization, and effect modification (Sections 10.9 and 10.10)
- The Mantel extension test (Section 10.10.4)
- Power and sample size estimation for stratified categorical data (Section 10.11)
- Greatly expanded section on assessing goodness of fit of regression models including residual analysis for both simple linear regression (Section 11.6) and multiple linear regression (Section 11.7.3)
- Partial regression coefficients (Section 11.7.1)
- Partial residual plots (Section 11.7.3)
- Relationship between t test methods, analysis of variance, analysis of covariance, and regression analysis (Section 11.8)
- Interval estimation for correlation coefficients (Section 11.11.3)
- Partial and multiple correlation (Section 11.12)
- Intraclass correlation coefficient (Section 11.13)
- Expanded discussion of multiple logistic regression including methods for prediction and assessment of goodness of fit (Sections 11.14.4 and 11.14.5)
- A new chapter on inference for person-time data (Chapter 13), including
- Measures of effect for person-time data (Section 13.1)
- Inference for stratified person-time data (Section 13.3)
- Power and sample size estimation for person-time data (Section 13.4)
- Testing for trend with incidence rate data (Section 13.5)
- Estimation of survival curves with the Kaplan-Meier estimator (Section 13.7)

The new sections and the expanded sections for this edition have been indicated by an asterisk in the Contents.

#### The Exercises

There are a total of 1600 exercises in the Fourth Edition (compared with 1300 in the Third Edition). Students have indicated that they would like to see more completely solved problems. As a result, 639 of the problems have been moved to a Study Guide to accompany the text given with complete solutions. 95 of these problems are given in a Miscellaneous Problems section and are randomly ordered so that they are not tied to a specific chapter in the book. This gives the student additional practice in determining "what method to use in what situation." The remaining 544 problems are related to specific chapters in the text. All problems have complete solutions. Approximately 900 problems remain in the text, including all data-set based problems. Brief solutions are given to 300 of these problems in the Answer section, and are indicated by an asterisk (\*) in the problem section of each chapter.

#### **A Method of Computation**

The method of handling computations in this edition of the book has also changed. All intermediate results are carried to full precision (10+ significant digits) even though they are presented with fewer significant digits (usually 2-3) in the text. Thus, intermediate results may seem to be inconsistent with final results in some instances, although this is not the case. This method allows for greater accuracy of final results and is the reason why there are slightly different results given for many calculations versus the previous editions, where intermediate results were carried to the same precision as shown in the text.

#### **Organization**

Fundamentals of Biostatistics, fourth edition, is organized as follows:

**Chapter 1** is an *introductory chapter* giving an outline of the development of an actual medical study I was involved with. It provides a unique sense of the role of biostatistics in the medical research process.

Chapter 2 concerns descriptive statistics and presents all the major numeric and graphic tools used for displaying medical data. This chapter is especially important for both consumers and producers of medical literature, since much of the actual communication of information is accomplished via descriptive material.

Chapters 3 through 5 discuss *probability*. The basic principles of probability are developed, and the most common probability distributions, such as the binomial and normal distributions, are introduced. These distributions are used extensively in the later chapters of the book.

Chapters 6 through 10 cover some of the basic methods of statistical inference.

**Chapter 6** introduces the concept of drawing random samples from populations. The difficult notion of a sampling distribution is also developed, including an introduction to the most common sampling distributions, such as the *t* and chi-square distributions. The basic methods of *estimation* are also presented, including an extensive discussion of confidence intervals.

**Chapters 7 and 8** contain the basic principles of *hypothesis testing*. The most elementary hypothesis tests for normally distributed data, such as the *t* test, are also fully discussed for one- and two-sample problems.

**Chapter 9** introduces the basic principles of the *analysis of variance* (ANOVA). The one-way analysis of variance fixed and random effects models are discussed as well as the analysis of data obtained using crossover designs.

Chapter 10 contains the basic concepts of *hypothesis testing* as applied to categorical data, including some of the most widely used statistical procedures, such as the chi-square test and Fisher's exact test.

**Chapter 11** develops the principles of *regression analysis*. The case of simple linear regression is thoroughly covered, and extensions are provided for the multiple regression case. Important sections on goodness-of-fit of regression models are also included. Multiple logistic regression is also discussed.

**Chapter 12** covers the basic principles of *nonparametric statistics*. The assumptions of normality are relaxed, and distribution-free analogues are developed for the tests in Chapters 7, 8, 9, and 11.

Chapter 13 introduces methods of analysis for person-time data. Included are methods for incidence rate data, as well as methods of survival analysis including the Kaplan-Meier survival curve estimator, the log rank test, and the Cox proportional hazards model.

The elements of study design are also discussed, including the concepts of matching, cohort studies, case-control studies, retrospective studies, prospective studies, and the sensitivity, specificity, and predictive value of screening tests. These designs are presented in the context of actual samples. In addition, specific sections on sample size estimation are provided for different statistical situations in Chapters 7, 8, 9, and 10.

A flowchart of appropriate methods of statistical inference on pages 671–675 provides an easy reference to the methods developed in this book. This flowchart is referred to at the end of each of Chapters 7 through 13 to give the student some perspective on how the methods in a particular chapter fit in with the overall collection of statistical methods introduced in this book.

In addition, an index summarizing all examples and problems used in this book is provided, grouped by *medical specialty*.

#### **Acknowledgments**

I am indebted to Debra Sheldon, Marie Sheehan, and Harry Taplin, who have been invaluable in helping to type this manuscript. I am indebted to those who reviewed the manuscript, among them: Stuart J. Anderson, University of Pittsburgh; Kenneth J. Koehler, Iowa State University; Donald J. Slymen, San Diego State University; and Craig D. Turnbull, University of North Carolina, Chapel Hill. I wish to thank Alex Kugushev, Jennie Burger, George Calmenson, and Linda Purrington, who were instrumental in providing editorial advice and in the preparation of the manuscript. I am indebted to my many colleagues at the Channing Laboratory, most notably Edward Kass, Frank Speizer, Charles Hennekens, Frank Polk, Ira Tager, Jerome Klein, James Taylor, Stephen Zinner, Scott Weiss, Frank Sacks, Walter Willett, Alvaro Munoz, Graham Colditz, and Susan Hankinson, and to my other colleagues at the Harvard Medical School, most notably Frederick Mosteller, Eliot Berson, Robert Ackerman, Mark Abelson, Arthur Garvey, Leo Chylack, Eugene Braunwald, and Arthur Dempster, who provided the inspiration for writing this book. Finally, I wish to acknowledge Leslie Miller, Andrea Wagner, Loren Fishman, and Roberta Shapiro, without whose clinical help the current edition of this book would not have been possible.

> Bernard Rosner Boston

### CONTENTS

#### CHAPTER 1 General Overview 1

■■ Reference, 4

#### CHAPTER 2 Descriptive Statistics 5

- 2.1 Introduction, 5
- 2.2 Measures of Central Location, 6
- 2.3 Some Properties of the Arithmetic Mean, 14
- 2.4 Measures of Spread, 15
- **2.5** Some Properties of the Variance and Standard Deviation, 21
- 2.6 The Coefficient of Variation, 23
- 2.7 Grouped Data, 24

- 2.8 Graphic Methods for Grouped Data, 29
- \*2.9 Case Study: Effects of Lead Exposure on Neurological and Psychological Function in Children, 35
- **2.10** Summary, 37
  - ■■ Problems, 38
  - ■■ References, 42

#### CHAPTER 3 Probability 43

- 3.1 Introduction, 43
- 3.2 Definition of Probability, 43
- 3.3 Some Useful Probabilistic Notation, 45
- \*3.4 The Multiplication Law of Probability, 47
- 3.5 The Addition Law of Probability, 49
- 3.6 Conditional Probability, 52

- 3.7 Bayes' Rule and Screening Tests, 56
- 3.8 Prevalence and Incidence, 61
- 3.9 Summary, 61
- ■■ Problems, 62
- ■■ References, 69

#### CHAPTER 4 Discrete Probability Distributions 71

- 4.1 Introduction, 71
- 4.2 Random Variables, 72
- **4.3** The Probability Mass Function for a Discrete Random Variable, 72
- **4.4** The Expected Value of a Discrete Random Variable, 74
- **4.5** The Variance of a Discrete Random Variable, 76
- **4.6** The Cumulative-Distribution Function of a Discrete Random Variable, 78
- 4.7 Permutations and Combinations, 79
- 4.8 The Binomial Distribution, 82

<sup>\*</sup>asterisks indicate new section or subsection for the Fourth Edition

#### VIII CONTENTS

- **4.9** Expected Value and Variance of the Binomial Distribution, 87
- 4.10 The Poisson Distribution, 88
- 4.11 Computation of Poisson Probabilities, 92
- **4.12** Expected Value and Variance of the Poisson Distribution, 94
- **4.13** Poisson Approximation to the Binomial Distribution, 95
- 4.14 Summary, 97
  - ■■ Problems, 97
  - ■■ References, 103

#### CHAPTER 5 Continuous Probability Distributions 105

- 5.4 Introduction, 105
- 5.2 General Concepts, 105
- 5.3 The Normal Distribution, 107
- **5.4** Properties of the Standard Normal Distribution, 110
- **5.5** Conversion from an  $N(\mu, \sigma^2)$  Distribution to an N(0, 1) Distribution, 115
- 5.6 Linear Combinations of Random Variables, 119

- **5.7** Normal Approximation to the Binomial Distribution, 121
- **5.8** Normal Approximation to the Poisson Distribution, 125
- 5.9 Summary, 130
- ■■ Problems, 133
- ■■ References, 140

#### CHAPTER 6 Estimation 141

- 6.1 Introduction, 141
- **6.2** The Relationship Between Population and Sample, 142
- 6.3 Random-Number Tables, 143
- \*6.4 Randomized Clinical Trials, 147
- 6.5 Estimation of the Mean of a Distribution, 151
- **6.6** Estimation of the Variance of a Distribution, 168

- **6.7** Estimation for the Binomial Distribution, 173
- \*6.8 Estimation for the Poisson Distribution, 178
- 6.9 One-Sided Confidence Intervals, 182
- **6.10** Summary, 184
  - ■■ Problems, 185
  - ■■ References, 190

#### CHAPTER 7 Hypothesis Testing: One-Sample Inference 191

- 7.1 Introduction, 191
- 7.2 General Concepts, 192
- 7.3 One-Sample Test for the Mean of a Normal Distribution with Known Variance: One-Sided Alternatives, 194
- 7.4 One-Sample Test for the Mean of a Normal Distribution with Known Variance: Two-Sided Alternatives, 203
- **7.5** One-Sample *t* Test, 207
- 7.6 The Power of a Test, 212
- 7.7 Sample-Size Determination, 219
- **7.8** The Relationship Between Hypothesis Testing and Confidence Intervals, 225
- **7.9** One-Sample  $\chi^2$  Test for the Variance of a Normal Distribution, 228

- **7.10** One-Sample Test for a Binomial Proportion, 231
- \***7.11** One-Sample Inference for the Poisson Distribution, 237
- 7.12 Summary, 243
  - ■■ Problems, 245
  - ■■ References, 250

#### CHAPTER 8 Hypothesis Testing: Two-Sample Inference 251

- 8.1 Introduction, 251
- 8.2 The Paired t Test, 253
- **8.3** Interval Estimation for the Comparison of Means from Two Paired Samples, 256
- **8.4** Two-Sample *t* Test for Independent Samples with Equal Variances, 257
- **8.5** Interval Estimation for the Comparison of Means from Two Independent Samples (Equal Variance Case), 261
- **8.6** Testing for the Equality of Two Variances, 263

- **8.7** Two-Sample *t* Test for Independent Samples with Unequal Variances, 270
- \*8.8 Case Study: Effects of Lead Exposure on Neurological and Psychological Function in Children, 276
- \*8.9 The Treatment of Outliers, 277
- **8.10** Estimation of Sample Size and Power for Comparing Two Means, 283
- 8.11 Summary, 285
  - ■■ Problems, 286
  - ■■ References, 297

#### CHAPTER 9 Multisample Inference 299

- **9.1** Introduction to the One-Way Analysis of Variance, 299
- 9.2 One-Way Analysis of Variance—Fixed-Effects Model, 299
- 9.3 Hypothesis Testing in One-Way ANOVA— Fixed-Effects Model, 301
- **9.4** Comparisons of Specific Groups in One-Way ANOVA, 306
- \*9.5 Case Study: Effects of Lead Exposure on Neurological and Psychological Function in Children, 319

- \*9.6 One-Way ANOVA—The Random-Effects Model, 322
- \*9.7 The Cross-Over Design, 329
- 9.8 Summary, 337
- ■■ Problems, 337
- ■■ References, 343

#### CHAPTER 10 Hypothesis Testing: Categorical Data 345

- **10.1** Introduction, 345
- **10.2** Two-Sample Test for Binomial Proportions, 346
- \*10.3 Study Design, 359

- \*10.4 Measures of Effect for Categorical Data, 361
  - 10.5 Fisher's Exact Test, 370
  - **10.6** Two-Sample Test for Binomial Proportions for Matched-Pair Data (McNemar's Test), 377

#### X CONTENTS

- **10.7** Estimation of Sample Size and Power for Comparing Two Binomial Proportions, 383
- **10.8**  $R \times C$  Contingency Tables, 392
- \*10.9 Confounding and Standardization, 399
- **10.10** Methods of Inference for Stratified Categorical Data—The Mantel-Haenszel Test, 404
- \*10.11 Power and Sample-Size Estimation for Stratified Categorical Data, 417

- 10.12 Chi-Square Goodness-of-Fit Test, 419
- 10.13 The Kappa Statistic, 423
- 10.14 Summary, 427
  - ■■ Problems, 428
  - ■■ References, 440

#### CHAPTER 11 Regression and Correlation Methods 443

- 11.1 Introduction, 443
- 11.2 General Concepts, 444
- **11.3** Fitting Regression Lines—The Method of Least Squares, 447
- **11.4** Inferences About Parameters from Regression Lines, 452
- **11.5** Interval Estimation for Linear Regression, 462
- \*11.6 Assessing the Goodness of Fit of Regression Lines, 466
- 11.7 Multiple Regression, 469
- \*11.8 Case Study: Effects of Lead Exposure on Neurological and Psychological Function in Children, 483

- 11.9 Two-Way Analysis of Variance, 496
- 11.10 The Correlation Coefficient, 503
- **11.11** Statistical Inference for Correlation Coefficients, 506
- \*11.12 Partial and Multiple Correlation, 516
- \*41.13 The Intraclass Correlation Coefficient, 517
- 11.14 Multiple Logistic Regression, 521
- **11.15** Summary, 540
  - ■■ Problems, 541
  - ■■ References, 549

#### CHAPTER 12 Nonparametric Methods 551

- 12.1 Introduction, 551
- 12.2 The Sign Test, 553
- 12.3 The Wilcoxon Signed-Rank Test, 558
- 12.4 The Wilcoxon Rank-Sum Test, 562
- 12.5 The Kruskal-Wallis Test, 569

- 12.6 Rank Correlation, 575
- **12.7** Summary, 579
- ■■ Problems, 580
- ■■ References, 584

#### CHAPTER 13 Hypothesis Testing: Person-Time Data 585

- \***13.1** Measures of Effect for Person-Time Data, 585
- **13.2** Two-Sample Inference for Incidence-Rate Data, 587
- \***13.3** Inference for Stratified Person-Time Data, 594
- \*13.4 Power and Sample-Size Estimation for Person-Time Data, 601

\*13.5 Testing for Trend-Incidence-Rate Data, 604

13.6 Introduction to Survival Analysis, 607

\***13.7** Estimation of Survival Curves: The Kaplan-Meier Estimator, 609

13.8 The Log-Rank Test, 614

13.9 The Proportional-Hazards Model, 621

13.10 Summary, 628

■■ Problems, 628

■■ References, 631

#### APPENDIX Tables 633

**1** Exact Binomial Probabilities  $Pr(X = k) = \binom{n}{k} p^k q^{n-k}$ , 635

**2** Exact Poisson Probabilities  $Pr(X = k) = \frac{e^{-\mu}\mu^k}{k!}, 640$ 

3 The Normal Distribution, 643

4 Table of 1000 Random Digits, 648

**5** Percentage Points of the *t* Distribution  $(t_{d,u})$ , 649

**6** Percentage Points of the Chi-Square Distribution  $(\chi_{d,u}^2)$ , 650

**7a** Exact Two-Sided 100%  $(1 - \alpha)$  Confidence Limits for Binomial Proportions  $(\alpha = .05)$ , 651

**7b** Exact Two-Sided 100%  $(1 - \alpha)$  Confidence Limits for Binomial Proportions  $(\alpha = .01)$ , 652

\*8 Confidence Limits for the Expectation of a Poisson Random Variable (μ), 653

**9** Percentage Points of the *F* Distribution  $(F_{d_1,d_2,p})$ , 654

\*10 Critical Values for the ESD (Extreme Studentized Deviate) Outlier Statistic (ESD<sub>1- $\alpha$ </sub>,  $\alpha = .05, .01$ ), 656

11 Fisher's z Transformation, 657

**12** Two-Tailed Critical Values for the Wilcoxon Signed-Rank Test, 657

**13** Two-Tailed Critical Values for the Wilcoxon Rank-Sum Test. 658

14 Critical Values for the Kruskal-Wallis Test Statistic (H) for Selected Sample Sizes for k = 3, 660

**15** Two-Tailed Upper Critical Values for the Spearman Rank-Correlation Coefficient  $(r_s)$ , 661

#### **Answers to Selected Problems 663**

#### Flowchart for Appropriate Methods of Statistical Inference 669

#### Index of Data Sets 676

#### Index 677

### **GENERAL OVERVIEW**

Statistics is the science whereby inferences are made about specific random phenomena on the basis of relatively limited sample material. The field of statistics can be subdivided into two main areas: mathematical statistics and applied statistics. Mathematical statistics concerns the development of new methods of statistical inference and requires detailed knowledge of abstract mathematics for its implementation. Applied statistics concerns the application of the methods of mathematical statistics to specific subject areas, such as economics, psychology, and public health. Biostatistics is the branch of applied statistics that concerns the application of statistical methods to medical and biological problems.

A good way to learn about biostatistics and its role in the research process is to follow the flow of a research study from its inception at the planning stage to its completion, which usually occurs when a manuscript reporting the results of the study is published. As an example, I will describe one such study in which I participated.

A friend called one morning and in the course of our conversation mentioned that he had recently used a new, automated blood-pressure device of the type seen in many banks, hotels, and department stores. The machine had read his average diastolic blood pressure on several occasions as 115 mm Hg; the highest reading was 130 mm Hg. I was horrified to hear of his experience, since if these readings were true, my friend might be in imminent danger of having a stroke or developing some other serious cardiovascular disease. I referred him to a clinical colleague of mine who, using a standard blood-pressure cuff, measured my friend's diastolic blood pressure as 90 mm Hg. The contrast in the readings aroused my interest, and I began to jot down the readings on the digital display every time I passed the machine at my local bank. I got the distinct impression that a large percentage of the reported readings were in the hypertensive range. Although one would expect that hypertensives would be more likely to use such a machine, I still believed that blood-pressure readings obtained with the machine might not be comparable with those obtained using standard methods of blood-pressure measurement. I spoke to Dr. B. Frank Polk about my suspicion and succeeded in interesting him in a small-scale evaluation of such machines. We decided to send a human observer who was well trained in blood-pressure measurement techniques to several of these machines. He would offer to pay subjects 50¢ for the cost of using the machine if they would agree to fill out a short questionnaire and have their blood pressure measured by both a human observer and the machine.

At this stage we had to make several important decisions, each of which would prove vital to the success of the study. The decisions were based on the following questions:

- (1) How many machines should we test?
- (2) How many people should we test at each machine?
- (3) In what order should the measurements be taken—should the human observer or the machine be used first? Ideally, we would have preferred to avoid this problem by taking both the human and machine readings simultaneously, but this procedure was logistically impossible.

- **(4)** What other data should we collect on the questionnaire that might influence the comparison between methods?
- (5) How should the data be recorded to facilitate their computerization at a later date?
- (6) How should the accuracy of the computerized data be checked?

We resolved these problems as follows:

- (1) and (2) We decided to test more than one machine (four to be exact), since we were not sure if the machines were comparable in quality. However, we wanted to sample enough subjects from each machine so that we would have an accurate comparison of the standard and automated methods for each machine. We tried to predict how large a discrepancy there might be between the two methods. Using the methods of sample-size determination discussed in this book, we calculated that we would need 100 subjects at each site to have an accurate comparison.
- (3) We then had to decide in what order the measurements should be taken for each person. According to some reports, one problem that occurs with repeated blood-pressure measurements is that people tense up at the initial measurement, yielding higher blood pressure than at subsequent repeated measurements. Thus, we would not always want to use the automated or manual method first, since the effect of the method would get confused with the order-of-measurement effect. A conventional technique that we used here was to **randomize** the order in which the measurements were taken, so that for any person it was equally likely that the machine or the human observer would take the first measurement. This random pattern could be implemented by flipping a coin or, more likely, by using a table of **random numbers** as appears in Table 4 of the Appendix.
- (4) We felt that the major extraneous factor that might influence the results would be body size, since we might have more difficulty getting accurate readings from people with fatter arms than from those with leaner arms. We also wanted to get some idea of the type of people who use these machines; so we asked questions about age, sex, and previous hypertensive history.
- (5) To record the data, we developed a coding form that could be filled out on site and from which data could be easily entered on a computer terminal for subsequent analysis. Each person in the study was assigned an identification (ID) number by which the computer could uniquely identify that person. The data on the coding forms were then keyed and verified. That is, the same form was entered twice, and a comparison was made between the two records to make sure they were the same. If the records were not the same, the form was reentered.
- (6) After data entry we ran some editing programs to ensure that the data were accurate. Checking each item on each form was impossible because of the large amount of data. Alternatively, we checked that the values for individual variables were within specified ranges and printed out aberrant values for manual checking. For example, we checked that all blood-pressure readings were at least 50 and no more than 300 and printed out all readings that fell outside this range.

After completing the data-collection, data-entry, and data-editing phases, we were ready to look at the results of the study. The first step in this process is to get a general feel for the data by summarizing the information in the form of several descriptive

statistics. This descriptive material can be numerical or graphical. If numerical, it can be in the form of a few summary statistics, which can be presented in tabular form or, alternatively, in the form of a **frequency distribution**, which lists each value in the data and how frequently it occurs. If graphical, the data are summarized pictorially and can be presented in one or more figures. The appropriate type of descriptive material will vary with the type of distribution considered. If the distribution is **continuous**, that is, if there are essentially an infinite number of possible values, as would be the case for blood pressure, then means and standard deviations might be the appropriate descriptive statistics. However, if the distribution is **discrete**, that is, if there are only a few possible values, as would be the case for sex, then percentages of people taking on each value would be the appropriate descriptive measure. In some cases both types of descriptive statistics are used for continuous distributions by condensing the range of possible values into a few groups and giving the percentage of people that fall into each group (e.g., the percentages of people that have blood pressures between 120 and 129 mm Hg and between 130 and 139 mm Hg, etc.).

In this study we decided first to look at mean blood pressure for each method at each of the four sites. Table 1.1 summarizes this information [1].

TABLE 1.1
Mean blood pressures
and differences
between machine and
human readings at
four locations

Location	Number of people	Systolic blood pressure (mm Hg)					
		Machine		Human		Difference	
		Mean	Standard deviation	Mean	Standard deviation	Mean	Standard deviation
A	98	142.5	21.0	142.0	18.1	0.5	11.2
В	84	134.1	22.5	133.6	23.2	0.5	12.1
C	98	147.9	20.3	133.9	18.3	14.0	11.7
D	62	135.4	16.7	128.5	19.0	6.9	13.6

Source: By permission of the American Heart Association, Inc.

You might notice from this table that we did not obtain meaningful data from all the 100 people interviewed at each site, since we could not obtain valid readings from the machine for many of the people. This type of missing-data problem is very common in biostatistics and should be anticipated at the planning stage when deciding on sample sizes (which was not done in this study).

Our next step in the study was to determine whether the apparent differences in blood pressure between machine and human measurements at two of the locations (C, D) were "real" in some sense or were "due to chance." This type of question falls into the area of **inferential statistics**. We realized that although there was a 14-mm Hg difference in mean systolic blood pressure between the two methods for the 98 people we interviewed at location C, this difference might not hold up if we interviewed 98 other people at a different time, and we wanted to have some idea as to the **error in the estimate** of 14 mm Hg. In statistical jargon this group of 98 people represents a **sample** from the **population** of all people who use that machine. We were interested in the population and we wished to use the sample to help us learn something about

the population. In particular, we wanted to know how different the estimated mean difference of 14 mm Hg in our sample was likely to be from the true mean difference in the population of all people who might use this machine. More specifically, we wanted to know if it was still possible that there was no underlying difference between the two methods and that our results were due to chance. The 14-mm Hg difference in our group of 98 people is referred to as an estimator of the true mean difference (d) in the population. The problem of inferring characteristics of a population from a sample is the central concern of statistical inference and is a major topic in this text. To accomplish this aim, we needed to develop a probability model, which would tell us how likely it is that we would obtain a 14-mm Hg difference between the two methods in a sample of 98 people if there were no real difference between the two methods over the entire population of users of the machine. If this probability were sufficiently small, then we would begin to believe that a real difference existed between the two methods. In this particular case, using a probability model based on the t distribution, we were able to conclude that this probability was less than 1 in 1000 for each of machines C and D. This probability was sufficiently small for us to conclude that there was a real difference between the automatic and manual methods of taking blood pressure for two of the four machines tested.

We used a statistical package to perform the preceding data analyses. A package is a collection of statistical programs that describe data and perform various statistical tests on the data. Currently the most widely used statistical packages include SAS, SPSS<sup>X</sup>, BMDP, and Minitab.

The final step in this study, after completing the data analysis, was to compile the results in the form of a publishable manuscript. Inevitably, because of space considerations, much of the material developed during the data-analysis phase was weeded out and only the essential items were presented for publication.

The review of this study should give you some idea of what medical research is about and what the role of biostatistics is in this process. The material in this text parallels the description of the data-analysis phase of the study described. Chapter 2 summarizes different types of descriptive statistics. In Chapters 3 through 5, some basic principles of probability and various probability models for use in later discussions of inferential statistics are presented. In Chapters 6 through 13, the major topics of inferential statistics as used in biomedical practice are discussed. Issues of study design or data collection are brought up only as they relate to other topics discussed in the text.

#### Reference

[1] Polk, B. F., Rosner, B., Feudo, R., & Vandenburgh, M. (1980). An evaluation of the Vita-Stat automatic blood pressure measuring device. *Hypertension*, 2(2), 221–227.

# DESCRIPTIVE STATISTICS

#### **SECTION 2.1** Introduction

The first step in looking at data is to describe the data at hand in some concise way. In smaller studies this step can be accomplished by listing each data point. In general, however, this procedure is tedious or impossible and, even if it were possible, would not give an overall picture of what the data look like.

#### EXAMPLE 2.1

**Cancer, Nutrition** Some investigators have proposed that consumption of vitamin A prevents cancer. To test this theory, a dietary questionnaire to collect data on vitamin-A consumption among 200 hospitalized cancer cases and 200 controls might be used. The controls would be matched on age and sex to the cancer cases and would be in the hospital at the same time for an unrelated disease. What should be done with these data after they are collected?

Before any formal attempt to answer this question can be made, the vitamin-A consumption among cases and controls must be described. Consider Figure 2.1. The **bar graphs** show visually that the controls have a higher vitamin-A consumption than the cases do, particularly in doses higher than the recommended daily allowance (RDA).

#### EXAMPLE 2.2

Pulmonary Disease Medical researchers have often suspected that passive smokers—people who themselves do not smoke but who live or work in an environment where others smoke—might have impaired pulmonary function as a result. In 1980 a research group in San Diego published results indicating that passive smokers did indeed have significantly lower pulmonary function than comparable nonsmokers who did not work in smoky environments [1]. As supporting evidence, the authors measured the carbon-monoxide (CO) concentrations in the working environments of passive smokers and of nonsmokers (where no smoking was permitted in the workplace) to see if the relative CO concentration changed over the course of the day. These results are displayed in the form of a scatter plot in Figure 2.2.

Figure 2.2 clearly shows that the CO concentrations in the two working environments are about the same early in the day but diverge widely in the middle of the day and then converge again after the working day is over at 7 P.M.

Graphic displays illustrate the important role of descriptive statistics, which is to quickly display data to give the researcher a clue as to the principal trends in the data and suggest hints as to where a more detailed look at the data, using the methods of inferential statistics, might be worthwhile. Descriptive statistics are also crucially important in conveying the final results of studies in written publications. Unless it is one of their primary interests, most readers will not have time to critically evaluate the work of others but will be influenced mainly by the descriptive statistics presented.