

# 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, cook-book 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

\*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

---

- 5.1 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
- \*11.13** 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 ( $\chi^2_{d,u}$ ), 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 ( $\mu$ ), 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 ( $ESD_{1-\alpha}$ ,  $\alpha = .05, .01$ ), 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**


---

