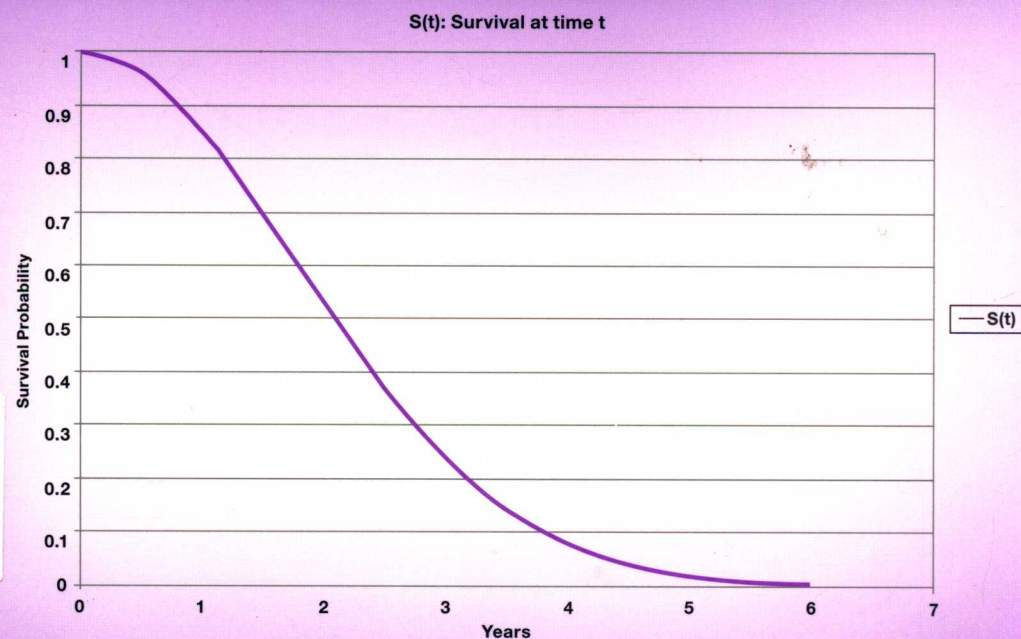


The Essentials of Biostatistics for Physicians, Nurses, and Clinicians

Michael R. Chernick



The Essentials of Biostatistics for Physicians, Nurses, and Clinicians

Michael R. Chernick

*Lankenau Institute for Medical Research
Wynnewood, PA*



 **WILEY**

A John Wiley & Sons, Inc., Publication

Copyright © 2011 by John Wiley & Sons, Inc. All rights reserved.

Published by John Wiley & Sons, Inc., Hoboken, New Jersey.

Published simultaneously in Canada.

No part of this publication may be reproduced, stored in a retrieval system, or transmitted in any form or by any means, electronic, mechanical, photocopying, recording, scanning, or otherwise, except as permitted under Section 107 or 108 of the 1976 United States Copyright Act, without either the prior written permission of the Publisher, or authorization through payment of the appropriate per-copy fee to the Copyright Clearance Center, Inc., 222 Rosewood Drive, Danvers, MA 01923, (978) 750-8400, fax (978) 750-4470, or on the web at www.copyright.com. Requests to the Publisher for permission should be addressed to the Permissions Department, John Wiley & Sons, Inc., 111 River Street, Hoboken, NJ 07030, (201) 748-6011, fax (201) 748-6008, or online at <http://www.wiley.com/go/permissions>.

Limit of Liability/Disclaimer of Warranty: While the publisher and author have used their best efforts in preparing this book, they make no representations or warranties with respect to the accuracy or completeness of the contents of this book and specifically disclaim any implied warranties of merchantability or fitness for a particular purpose. No warranty may be created or extended by sales representatives or written sales materials. The advice and strategies contained herein may not be suitable for your situation. You should consult with a professional where appropriate. Neither the publisher nor author shall be liable for any loss of profit or any other commercial damages, including but not limited to special, incidental, consequential, or other damages.

For general information on our other products and services or for technical support, please contact our Customer Care Department within the United States at (800) 762-2974, outside the United States at (317) 572-3993 or fax (317) 572-4002.

Wiley also publishes its books in a variety of electronic formats. Some content that appears in print may not be available in electronic formats. For more information about Wiley products, visit our web site at www.wiley.com.

Library of Congress Cataloging-in-Publication Data:

Chernick, Michael R.

The essentials of biostatistics for physicians, nurses, and clinicians / Michael R. Chernick.

p. ; cm.

Includes bibliographical references and index.

ISBN 978-0-470-64185-9 (pbk.)

1. Biometry. I. Title.

[DNLM: 1. Biostatistics. WA 950]

QH323.5C484 2011

570.1'5195-dc22

2011002198

oBook ISBN: 978-1-118-07195-3

ePDF ISBN: 978-1-118-07193-9

ePub ISBN: 978-1-118-07194-6

Printed in Singapore.

10 9 8 7 6 5 4 3 2 1



The Essentials of Biostatistics for Physicians, Nurses, and Clinicians

Preface

I have taught biostatistics in the health sciences and published a book in 2003 with Wiley on that topic. That book is a textbook for upper-level undergraduates and graduate students in the health science departments at universities. Since coming to the Lankenau Institute 17 months ago, I was tasked to prepare a course in biostatistics for nurses and physicians (particularly the hospital residents and fellows that do medical research). I quickly learned that although the material in my book was relevant, it contained too much material and was not in a digestible form for them. I prepared a six-lecture course (1 hour each) for physicians, and a two-lecture course for the nurses. To prevent boredom, I introduced some funny but educational cartoon slides. The course currently exists and has been refined as PowerPoint presentations and has been moderately successful. I also am starting a similar course at statistics.com.

The physicians and nurses have a busy schedule, and what they need is a concise and clearly explained set of lectures that cover only the areas of statistics that are essential to know about in medical research. This means topics that are not taught in traditional introductory statistics courses. So Kaplan–Meier curves, repeated measures analysis of variance, hazard ratios, contingency tables, logrank tests, bioequivalence, cross-over designs, noninferiority, selection bias, and group sequential methods are all included, but they are introduced on a conceptual level without the need for theory. It is when and why these methods work that they need to know, and not a detailed account of how they work mathematically. I feel that it would be appropriate to have a textbook for such a course that can be taught in-house at research centers or online courses. The book is intended to be approximately 160 pages along with suitable references.

I am very grateful to Professor Marlene Egger, who carefully reviewed the manuscript and made several wonderful suggestions that helped with the clarity and improved the content of the book.

Michael R. Chernick

Contents

Preface

ix

1. The What, Why, and How of Biostatistics in Medical Research

1

- 1.1 Definition of Statistics and Biostatistics, 1
- 1.2 Why Study Statistics?, 3
- 1.3 The Medical Literature, 9
- 1.4 Medical Research Studies, 11
 - 1.4.1 Cross-sectional studies including surveys, 11
 - 1.4.2 Retrospective studies, 12
 - 1.4.3 Prospective studies other than clinical trials, 12
 - 1.4.4 Controlled clinical trials, 12
 - 1.4.5 Conclusions, 13
- 1.5 Exercises, 14

2. Sampling from Populations

15

- 2.1 Definitions of Populations and Samples, 17
- 2.2 Simple Random Sampling, 18
- 2.3 Selecting Simple Random Samples, 19
- 2.4 Other Sampling Methods, 27
- 2.5 Generating Bootstrap Samples, 28
- 2.6 Exercises, 32

3. Graphics and Summary Statistics

34

- 3.1 Continuous and Discrete Data, 34
- 3.2 Categorical Data, 35
- 3.3 Frequency Histograms, 35
- 3.4 Stem-and-Leaf Diagrams, 38
- 3.5 Box Plots, 39
- 3.6 Bar and Pie Charts, 39
- 3.7 Measures of the Center of a Distribution, 42

3.8 Measures of Dispersion, 46

3.9 Exercises, 50

4. Normal Distribution and Related Properties **51**

4.1 Averages and the Central Limit Theorem, 51

4.2 Standard Error of the Mean, 53

4.3 Student's t -Distribution, 53

4.4 Exercises, 55

5. Estimating Means and Proportions **58**

5.1 The Binomial and Poisson Distributions, 58

5.2 Point Estimates, 59

5.3 Confidence Intervals, 62

5.4 Sample Size Determination, 65

5.5 Bootstrap Principle and Bootstrap Confidence Intervals, 66

5.6 Exercises, 69

6. Hypothesis Testing **72**

6.1 Type I and Type II Errors, 73

6.2 One-Tailed and Two-Tailed Tests, 74

6.3 P -Values, 74

6.4 Comparing Means from Two Independent Samples:
Two-Sample t -Test, 75

6.5 Paired t -Test, 76

6.6 Testing a Single Binomial Proportion, 78

6.7 Relationship Between Confidence Intervals and Hypothesis Tests, 79

6.8 Sample Size Determination, 80

6.9 Bootstrap Tests, 81

6.10 Medical Diagnosis: Sensitivity and Specificity, 82

6.11 Special Tests in Clinical Research, 83

6.11.1 Superiority tests, 84

6.11.2 Equivalence and bioequivalence, 84

6.11.3 Noninferiority tests, 86

6.12 Repeated Measures Analysis of Variance and Longitudinal Data Analysis, 86

6.13 Meta-Analysis, 88

6.14 Exercises, 92

7. Correlation, Regression, and Logistic Regression	95
7.1 Relationship Between Two Variables and the Scatter Plot, 96	
7.2 Pearson's Correlation, 99	
7.3 Simple Linear Regression and Least Squares Estimation, 101	
7.4 Sensitivity to Outliers and Robust Regression, 104	
7.5 Multiple Regression, 111	
7.6 Logistic Regression, 117	
7.7 Exercises, 122	
8. Contingency Tables	127
8.1 2×2 Tables and Chi-Square, 127	
8.2 Simpson's Paradox in the 2×2 Table, 129	
8.3 The General $R \times C$ Table, 132	
8.4 Fisher's Exact Test, 133	
8.5 Correlated Proportions and McNemar's Test, 136	
8.6 Relative Risk and Odds Ratio, 138	
8.7 Exercises, 141	
9. Nonparametric Methods	145
9.1 Ranking Data, 146	
9.2 Wilcoxon Rank-Sum Test, 146	
9.3 Sign Test, 149	
9.4 Spearman's Rank-Order Correlation Coefficient, 150	
9.5 Insensitivity of Rank Tests to Outliers, 153	
9.6 Exercises, 154	
10. Survival Analysis	158
10.1 Time-to-Event Data and Right Censoring, 159	
10.2 Life Tables, 160	
10.3 Kaplan–Meier Curves, 164	
10.3.1 The Kaplan–Meier curve: a nonparametric estimate of survival, 164	
10.3.2 Confidence intervals for the Kaplan–Meier estimate, 165	
10.3.3 The logrank and chi-square tests: comparing two or more survival curves, 166	

- 10.4 Parametric Survival Curves, 168
 - 10.4.1 Negative exponential survival distributions, 168
 - 10.4.2 Weibull family of survival distributions, 169
- 10.5 Cox Proportional Hazard Models, 170
- 10.6 Cure Rate Models, 171
- 10.7 Exercises, 173

Solutions to Selected Exercises	175
Appendix: Statistical Tables	192
References	204
Author Index	209
Subject Index	211

The What, Why, and How of Biostatistics in Medical Research

1.1 DEFINITION OF STATISTICS AND BIOSTATISTICS

The *Oxford Dictionary of Statistics* (2002, p. 349) defines statistics as “The science of collecting, displaying, and analyzing data.” Statistics is important in any scientific endeavor. It also has a place in the hearts of fans of sports, particularly baseball. Roger Angel in his baseball book, *Late Innings*, says “Statistics are the food of love.”

Biostatistics is the branch of statistics that deals with biology, both experiments on plants, animals, and living cells, and controlled experiments on humans, called clinical trials. Statistics is classified by scientific discipline because in addition to many standard methods that are common to statistical problems in many fields, special methods have been developed primarily for certain disciplines. So to illustrate, in biostatistics, we study longitudinal data, missing data models, multiple testing, equivalence and noninferiority testing, relative risk and odds ratios, group sequential and adaptive designs, and survival analysis, because these types of data and methods arise in clinical trials and other medical studies. Engineering statistics considers tolerance intervals and design of experiments. Environmental statistics has a concentration in

The Essentials of Biostatistics for Physicians, Nurses, and Clinicians,
First Edition. Michael R. Chernick.

© 2011 John Wiley & Sons, Inc. Published 2011 by John Wiley & Sons, Inc.

the analysis of spatial data, and so does geostatistics. Econometrics is the branch of statistics studied by economists, and deals a lot with forecasting and time series.

Statisticians are professionals trained in the collection, display, and analysis of data and the distribution theory that characterizes the variability of data. To become a good applied statistician, one needs to learn probability theory and the methods of statistical inference as developed by Sir Ronald A. Fisher, Jerzy Neyman, Sir Harold Jeffreys, Jimmie Savage, Bruno deFinetti, Harald Cramer, Will Feller, A. N. Kolmogorov, David Blackwell, Erich Lehmann, C. R. Rao, Karl and Egon Pearson, Abraham Wald, George Box, William Cochran, Fred Mosteller, Herman Chernoff, David Cox, and John Tukey in the twentieth century. These are some of the major developers of the foundations of probability and statistics. Of course, when selecting a list of famous contributors like this, many have been unintentionally omitted. In the late twentieth century and early twenty-first century, computer-intensive statistics arose, and a partial list of the leaders of that development are Brad Efron, Leo Brieman, David Freedman, Terry Speed, Jerry Friedman, David Siegmund, and T. L. Lai. In the area of biostatistics, we should mention Thomas Fleming, Stuart Pocock, Nathan Mantel, Peter Armitage, Shein-Chung Chow, Jen-pei Liu, and Gordon Lan. You will be introduced to these and other famous probabilists and statisticians in this book. An applied statistician must also become familiar with at least one scientific discipline in order to effectively consult with scientists in that field.

Statistics is its own discipline because it is much more than just a set of tools to analyze data. Although statistics requires the tools of probability, which are mathematical, it should not be thought of as a branch of mathematics. It is the appropriate way to summarize and analyze data when the data contains an element of uncertainty. This is very common when measurements are taken, since there is a degree of inaccuracy in every measurement. Statisticians develop mathematical models to describe the phenomena being studied. These models may describe such things as the time a bus will arrival at a scheduled stop, how long a person waits in line at a bank, the time until a patient dies or has a recurrence of a disease, or future prices of stocks, bonds, or gasoline.

Based on these models, the statistician develops methods of estimation or tests of hypotheses to solve certain problems related to the data.

Because almost every experiment involves uncertainty, statistics is the scientific method for quantitative data analysis.

Yet in the public eye, statistics and statisticians do not have a great reputation. In the course of a college education, students in the health sciences, business, psychology, and sociology are all required to take an introductory statistics course. The comments most common from these students are “this is the most boring class I ever took” and “it was so difficult, that I couldn’t understand any of it.” This is the fault of the way the courses are taught and not the fault of the subject. An introductory statistics course can be much easier to understand and more useful to the student than, say, a course in abstract algebra, topology, and maybe even introductory calculus. Yet many people don’t view it that way.

Also, those not well trained in statistics may see articles in medicine that are contradictory but still make their case through the use of statistics. This causes many of us to say “You can prove anything with statistics.” Also, there is that famous quote attributed to Disraeli. “There are lies, damn lies and statistics.” In 1954, Darrell Huff wrote his still popular book, *How to Lie with Statistics*. Although the book shows how graphs and other methods can be used to distort the truth or twist it, the main point of the book is to get a better understanding of these methods so as not to be fooled by those who misuse them. Statisticians applying valid statistical methods will reach consistent conclusions. The data doesn’t lie. It is the people that manipulate the data that lie. Four books that provide valuable lessons about misusing statistics are Huff (1954), Campbell (1974), Best (2001), and Hand (2008).

1.2 WHY STUDY STATISTICS?

The question is really why should medical students, physicians, nurses, and clinicians study statistics? Our focus is on biostatistics and the students we want to introduce it to. One good reason to study statistics is to gain knowledge from data and use it appropriately. Another is to make sure that we are not to be fooled by the lies, distortions, and misuses in the media and even some medical journals. The medical journals now commonly require good statistical methods as part of a research paper, and the sophistication of the methods used is greater.

So we learn statistics so that we know what makes sense when reading the medical literature, and in order to publish good research.

We also learn statistics so that we can provide intelligent answers to basic questions of a statistical nature. For many physicians and nurses, there is a fear of statistics. Perhaps this comes from hearing horror stories about statistics classes. It also may be that you have seen applications of statistics but did not understand it because you have no training. So this text is designed to help you conquer your fear of statistics. As you learn and gain confidence, you will see that it is logical and makes sense, and is not as hard as you first thought.

Major employers of statisticians are the pharmaceutical, biotechnology, and medical device companies. This is because the marketing of new drugs, biologics, and most medical devices must be approved by the U.S. Food and Drug Administration (FDA), and the FDA requires the manufacturers to demonstrate through the use of animal studies and controlled clinical trials the safety and effectiveness of their product. These studies must be conducted using valid statistical methods. So any medical investigator involved in clinical trials sponsored by one of these companies really needs to understand the design of the trial and the statistical implications of the design and the sample size requirements (i.e., number of patients need in the clinical trial). This requires at least one basic biostatistics course or good on-the-job training.

Because of uncontrolled variability in any experimental situation, statistics is necessary to organize the data and summarize it in a way so that signals (important phenomena) can be detected when corrupted by noise. Consequently, bench scientists as well as clinical researchers need some acquaintance with statistics. Most medical discoveries need to be demonstrated using statistical hypothesis testing or confidence interval estimation. This has increased in importance in the medical journals. Simple t -tests are not always appropriate. Analyses are getting much more sophisticated. Death and other time-to-event data require statistical survival analysis methods for comparison purposes.

Most scientific research requires statistical analysis. When Dr. Riffenburgh (author of the text *Statistics in Medicine*, 1999) is told by a physician “I’m too busy treating patients to do research,” he answers, “When you treat a patient, you have treated a patient. When you do research, you have treated ten thousand patients.”

In order to amplify these points, I will now provide five examples from my own experience in the medical device and pharmaceutical

industries where a little knowledge of statistics would have made life easier for some of my coworkers.

In the first scenario, suppose you are the coordinator for a clinical trial on an ablation catheter. You are enrolling subjects at five sites. You want to add a new site to help speed up enrollment. The IRB for the new site must review and approve your protocol for the site to enter your study. A member of the IRB asks what stopping rule you use for safety. How do you respond? You don't even know what a stopping rule is or even that the question is related to statistics! By taking this course, you will learn that statisticians construct stopping rules based upon accumulated data. In this case, there may be safety issues, and the stopping rule could be based on reaching a high number of adverse events. You won't know all the details of the rule or why the statistician chose, it but you will at least know that the statistician is the person who should prepare the response for the IRB.

Our second example involves you as a regulatory affairs associate at a medical device company that just completed an ablation trial for a new catheter. You have submitted your premarket approval application (PMA). In the statistical section of the PMA, the statistician has provided statistical analysis regarding the safety and efficacy of your catheter in comparison to other marketed catheters. A reviewer at the FDA sent you a letter asking why Peto's method was not used instead of Greenwood's approximation. You do not know what these two methods are or how they apply.

From this course, you will learn about survival analysis. In studying the effectiveness of an ablation procedure, we not only want to know that the procedure stopped the arrhythmia (possibly atrial fibrillation), but also that the arrhythmia does not recur. Time to recurrence is one measure of efficacy for the treatment. Based on the recurrence data from the trial, your statistician constructs a time-to-event curve called the Kaplan–Meier curve.

If we are interested in the probability of recurrence within 1 year, then the Peto and Greenwood methods are two ways to get approximate confidence intervals for it. Statistical research has shown differences in the properties of these two methods for obtaining approximate confidence intervals for survival probabilities. As an example, Greenwood's estimate of the lower confidence bound can be too high in situations where the number of subjects still at risk at the time point of interest is small.

374 Analysis of survival times

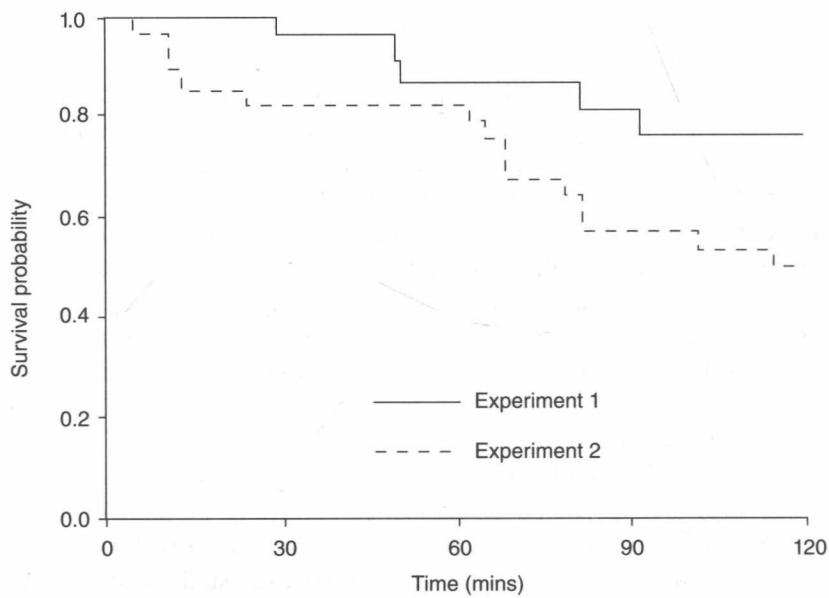


Figure 1.1. Example of a Kaplan–Meier curve. Taken from Altman (1991), *Practical Statistics for Medical Research*. Chapman and Hall/CRC, p. 374.

In these situations, Peto’s method gives a better estimate of this lower bound. In general, neither method is always superior to the other. Since the FDA posed this question, the statistician would opt to provide the Peto estimate in addition to Greenwood for the FDA to compare the two lower confidence bounds. Knowing these simple facts would help you deal with the FDA question quickly, effectively, and accurately (Fig. 1.1).

In situation 3, you are in regulatory affairs and are reviewing an FDA letter about a PMA submission. The FDA wants you to report results in terms of confidence intervals, in addition to the p -values, before they give final approval to the treatment. You recognize this as a statistical question, but are worried because if it takes significant time to supply the request, the launch date of the new device will be delayed and will upset marketing’s plans. You don’t even know what a confidence interval is!

In this case, since you have the necessary data to do the binomial test on success probability, you can easily compute an exact confidence

interval. Your statistician can provide this for you in less than 1 day and you are greatly relieved.

In situation 4, you are a clinical research associate in the middle of an important phase III trial. Based upon a data analysis done by the statistics group and an agreement with the FDA prior to the trial, the primary endpoint can be changed from a condition at the 6-month follow-up visit to that same condition at the 3-month follow-up visit. This is great news, because it means that the trial can be finished sooner!

There is a problem though. The protocol only required follow-up visits at 2 weeks and 6 months, and the 3-month follow-up was optional. Unfortunately, some sites opted not to conduct the 3-month follow-up. Your clinical manager now wants you to have all the patients that are past the 3-month time point since the procedure was done and did not have the 3-month follow-up to come in for an unscheduled visit. When you requested that the investigators do this, a nurse and one investigator balked at the idea and demanded to know why this is necessary. You need an answer from your statistician!

To placate the investigator, the statistician tells the investigator that they could not use the 3-month follow-up initially because the FDA had not seen data to indicate that a 3-month follow-up would be enough to determine long-term survival. However, during the early part of the trial, the statistician was able to find relevant survival curves to indicate the survival probability flattens out at 3 months' duration. This was enough to convince the FDA that the 3-month endpoint was sufficient to determine long-term survival. If we now have the unscheduled visits, these could be the subjects' last visit, and many subjects will not need a 6-month follow-up, allowing a shorter accrual time and a chance to get the product to market faster.

This explanation helped, but the problem could have been avoided had the clinician had the foresight to see the importance of making the 3-month follow-up mandatory in the protocol. The investigator was pleased because although it would cost more to add these unscheduled visits, this would be more than compensated by the dropping of the 6-month follow-up, for those getting the unscheduled visit, and possibly some others.

In the last situation (situation 5), imagine you are the VP of the Clinical and Regulatory Affairs Departments at a medical device company. Your company hired a contract research organization (CRO) to run a blinded randomized control phase III clinical trial. You have a