

Robert Hirsch

Analysis of Epidemiologic Data Using R

Synthesis Lectures on Mathematics & Statistics

Series Editor

Steven G. Krantz, Department of Mathematics, Washington University, Saint Louis, MO, USA

This series includes titles in applied mathematics and statistics for cross-disciplinary STEM professionals, educators, researchers, and students. The series focuses on new and traditional techniques to develop mathematical knowledge and skills, an understanding of core mathematical reasoning, and the ability to utilize data in specific applications.

Robert Hirsch

Analysis of Epidemiologic Data Using R

Robert Hirsch
Overland Park, KS, USA

ISSN 1938-1743 ISSN 1938-1751 (electronic)
Synthesis Lectures on Mathematics & Statistics
ISBN 978-3-031-41913-3 ISBN 978-3-031-41914-0 (eBook)
<https://doi.org/10.1007/978-3-031-41914-0>

© The Editor(s) (if applicable) and The Author(s), under exclusive license to Springer Nature Switzerland AG 2024

This work is subject to copyright. All rights are solely and exclusively licensed by the Publisher, whether the whole or part of the material is concerned, specifically the rights of translation, reprinting, reuse of illustrations, recitation, broadcasting, reproduction on microfilms or in any other physical way, and transmission or information storage and retrieval, electronic adaptation, computer software, or by similar or dissimilar methodology now known or hereafter developed.

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

The publisher, the authors, and the editors are safe to assume that the advice and information in this book are believed to be true and accurate at the date of publication. Neither the publisher nor the authors or the editors give a warranty, expressed or implied, with respect to the material contained herein or for any errors or omissions that may have been made. The publisher remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.

This Springer imprint is published by the registered company Springer Nature Switzerland AG
The registered company address is: Gewerbestrasse 11, 6330 Cham, Switzerland

Paper in this product is recyclable.

Preface

I had two goals in writing this text. One was to bring together methods of analysis for epidemiologic data in a way that did not require mathematics to understand the principles. This was accomplished with the only math being the methods themselves. The second goal was to make these methods accessible using R. This also was accomplished limited only by the existence of programs.

In addition to chapters addressing the statistical methods, there are four chapters addressing related topics. First, I assume that the reader has no familiarity with R, so there is a chapter that describes how to access R and how to do the types of things required to perform an analysis of epidemiologic data. Next, I assume that the reader has little familiarity with statistical methods, so there is a chapter describing the statistical logic behind estimation and hypothesis testing. Next, I assume that the reader has little familiarity with causal inference, so there is a chapter describing the concept of causation and how statistical methods can be used to support a causal conclusion. Finally, I assume that the reader has little familiarity with epidemiologic research designs, so there is a chapter describing the most common experimental and observational studies.

Where my assumptions about the reader are incorrect, the chapters that discuss topics with which the reader is familiar can be skipped without sacrificing appreciation of the remaining five chapters. In essence, the experienced reader can start with Chap. 5.

In selecting the topics of the last five chapters, I tried to cover all of the most common methods to analyze epidemiologic data. These chapters begin by discussing measures of disease frequency and measures of association using 2×2 tables as the structure. Then, I discuss life table analysis as a method for interpreting incidence data as risks. Next, I address stratified analysis beginning with an introduction to confounding. Finally, I present regression analysis as another method to control for confounding.

Overland Park, USA

Robert Hirsch

Acknowledgments I appreciate my clients who have shared their frustration with trying to understand statistical methods for epidemiologic data and their desire to be able to perform some statistical procedures without having to involve a statistician.

Notices

The examples in this book are based on fictitious data and should not be taken as a reflection of real relationships. Those data have been created only to illustrate statistical principles.

Contents

1	Introduction to R	1
1.1	Getting R on Your Computer	2
1.2	Using RStudio	2
1.3	R Datasets	3
1.4	Graphics	6
1.4.1	Scatterplots	6
1.4.2	Legends	9
1.4.3	Step Function Graphs	11
2	Overview of Statistical Logic	13
2.1	Estimation	14
2.2	Inference	15
3	Causal Inference	17
3.1	Sufficient Causes	18
3.2	Detection of a Causal Relationship	19
3.2.1	Bradford Hill's Criteria for Causation	20
3.2.2	Strength of Association	20
3.2.3	Consistency	20
3.2.4	Specificity	21
3.2.5	Temporal Sequence	21
3.2.6	Biologic Gradient	21
3.2.7	Biologic Rationale	21
3.2.8	Coherence	22
3.2.9	Experimental Evidence	22
3.2.10	Analogous Evidence	22
4	Design of Epidemiologic Studies	23
4.1	Experimental Studies	24
4.1.1	Clinical Trials	25
4.1.2	Field Trials	25

4.1.3	Community Intervention Trials	25
4.1.4	Cluster Randomized Trials	26
4.2	Observational Studies	26
4.2.1	Cohort Studies	26
4.2.2	Case–Control Studies	27
4.2.3	Cross-Sectional Studies	28
5	Measures of Disease Frequency	29
5.1	Prevalence	30
5.2	Risk	32
5.3	Incidence	34
6	2 × 2 Tables	39
6.1	Constructing 2 × 2 Tables	40
6.2	Inference	43
6.3	Interval Estimation	45
6.4	Paired Data	47
6.5	Incidence	50
7	Life Tables	57
7.1	Construction of Life Tables	58
7.2	Point Estimation	59
7.3	Inference	61
7.4	Interval Estimation	63
7.5	Survival Plots	64
7.6	Cutler-Ederer Tables	65
8	Stratified Analysis	69
8.1	Confounding	70
8.2	Summary Estimates	72
8.3	Hypothesis Test	76
8.4	Interval Estimates	77
8.5	Effect Modification	78
9	Regression Analysis	83
9.1	Logistic Regression	84
9.2	Cox Proportional Hazards Regression	89
	Answers to Exercises	93
	Index	107



Introduction to R

1

Abstract

R is free statistical software that is available for download on the web. It is usually run from an interface called “RStudio.” RStudio is set up with four quadrants. The lower left is where the programs are run. The lower right is where the documentation is available and where graphic output appears. The upper right is where datasets are listed, and the upper left is where you can see the content of datasets. Creation of datasets usually starts with the creation of vectors in the lower left quadrant. These vectors can be combined to create a dataset or data frame. A special kind of dataset is stored in a 2×2 matrix. These are often of use in epidemiology. We will be using several programs. These are arranged in packages. The “stats” package is the one we will use most often. Click on the name in the lower right quadrant to see a list of the programs available. Clicking on the name of the program invokes documentation on that program. Graphics are requested in the lower left quadrant, and they appear in the lower right quadrant. We use the “plot” command to create graphs. The plot command has several options used to customize your graph. The graphs can have additional data displayed using the “points” command. You can create a legend for your graph using the “legend” command. A special kind of graph used in epidemiology can be created using the “stepfun” function.

R is a collection of programs for statistical analysis and data management. It was first developed by two New Zealand statisticians, Ross Ihaka and Robert Gentleman and released in 1993.¹ It has since been expanded and modified by the R Core Team and the

¹ R got its name from the first letter of the developers’ first names.