# Epidemiologic Research

## PRINCIPLES AND QUANTITATIVE METHODS

Kleinbaum
Kupper
Morgenstern

# Epidemiologic Research

## PRINCIPLES AND QUANTITATIVE METHODS

David G. Kleinbaum

Lawrence L. Kupper

Hal Morgenstern

# Epidemiologic Research

# Table List

# Preface

This text discusses the principles and methods involved in the planning, analysis, and interpretation of epidemiologic research studies. Our purpose in writing this book is to provide the applied researcher with a synthesis of current methodological thought and practice. In particular, we focus our discussion on quantitative (including statistical) issues that arise during the course of an epidemiologic investigation. Instead of restricting our focus just to statistical techniques and their various applications, we also address issues of study *design, measurement,* and *validity*. Where appropriate, we describe statistical considerations relevant to each such issue.

**AIM**

Throughout the text, we emphasize that validity should be the primary goal in an epidemiologic study, even if this means sacrificing generalizability. A study free of bias, even if restricted in scope, is preferred to a more general study with unresolvable validity problems. Sophisticated statistical analyses mean nothing if the data are unreliable. Or, to put it another way, "good statistical analyses do not salvage poor data."

For several of the topics we discuss (e.g., the treatment of confounding, the use of matching, multivariable analyses), there are differing points of view in the literature. Even when there is reasonable consensus, there are rarely clear-cut recommendations that can be made. Thus, rather than try to find the nonexistent "ideal" solution, we provide a quantitative framework through which such issues can be investigated. Whenever possible, we discuss both the advantages and disadvantages of various approaches to a given methodological problem.

**PHILOSOPHY AND APPROACH**

**AUDIENCE AND PRE-REQUISITES**

This book was written with several audiences in mind. Our primary audience includes epidemiologists and other health professionals, as well as graduate students in the health sciences. We assume that this audience has a good understanding of the basic principles of epidemiology. A second audience includes researchers or students in other human science disciplines. Finally, a third audience consists of statisticians and biostatisticians interested in the application of statistics to epidemiologic research.

The degree of mathematical sophistication required, particularly with regard to statistical prerequisites, will depend on the reader's learning objectives for each chapter. We have written primarily at the level of persons with a good understanding of basic statistical procedures of data analysis typically covered in a two-course sequence in applied biostatistics. It would also be advantageous—through not essential—for the reader to have a basic knowledge of the biostatistical methods used to analyze epidemiologic data. Because we do not draw heavily on statistical inference principles prior to Part III, the reader with minimal statistical knowledge will be able to follow the discussion of basic methodological principles given in Parts I and II.

**OUTLINE OF TEXT**

This book is divided into three parts, plus an introductory chapter (Chapter 1). This overview chapter provides a nonquantitative discussion of key methodological issues that are addressed in the remainder of the text.

Part I (Chapters 2–9) contains the epidemiologic foundation for the more quantitatively oriented (including statistical) treatment of methodological issues addressed in Parts II and III. Of the eight chapters in Part I, only Chapters 6, 7, and 8 contain material (including notation) that is absolutely prerequisite for Parts II and III.

Part II (Chapters 10–14) provides a quantitative, conceptual framework for evaluating validity. Part III (Chapters 15–24) provides a detailed quantitative discussion of procedures and strategies for the design and analysis of epidemiologic research studies.

We strongly recommend that the reader work the exercises provided at the end of the chapters. These exercises are intended to help the reader understand and apply the principles and methods presented in the text. Abbreviated answers to selected problems are provided in Appendix A. A detailed solutions manual for the exercises is available from the publisher.

**ACKNOWL-EDGMENTS**

We wish to acknowledge several people who contributed to the preparation of this text. The authors developed much of their interest in quantitative epidemiology from exposure to the pioneering and sometimes controversial research contributions of Olli Miettinen at Harvard University. They have also gained insight from the work of Norman Breslow and Nicholas Day, as well as from many others too numerous to list individually. We thank Beverly Cole for typing the many drafts and revisions of all chapters, for her assistance in helping us to communicate efficiently

*To our parents*

Joslyn and Janet Kleinbaum

Louis and Sylvia Kupper

Joseph and Edith Morgenstern

# Contents