# Applied Missing Data Analysis in the Health Sciences
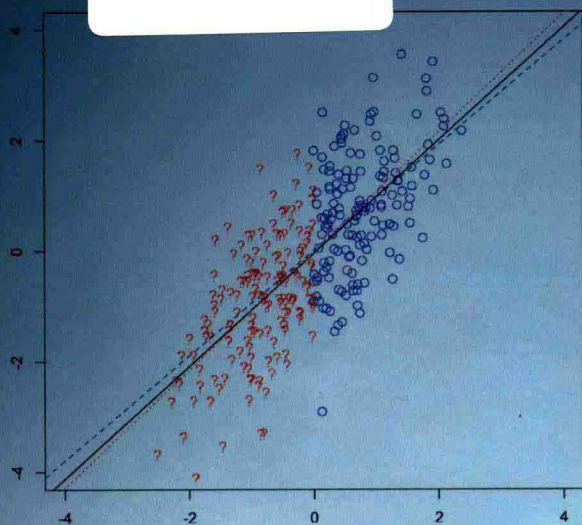


**XIAO-HUA ZHOU**
**CHUAN ZHOU**
**DANPING LIU**
**XIAOBO DING**

WILEY

# APPLIED MISSING DATA ANALYSIS IN HEALTH SCIENCES

**XIAO-HUA ZHOU**
University of Washington

**CHUAN ZHOU**
University of Washington

**DANPING LIU**
National Institutes of Health

**XIAOBO DING**
Chinese Academy of Sciences

WILEY

# APPLIED MISSING DATA ANALYSIS IN HEALTH SCIENCES

# WILEY SERIES IN STATISTICS IN PRACTICE

Advisory Editor, MARIAN SCOTT, *University of Glasgow, Scotland, UK*

Founding Editor, VIC BARNETT, *Nottingham Trent University, UK*

*Statistics in Practice* is an important international series of texts which provide detailed coverage of statistical concepts, methods, and worked case studies in specific fields of investigation and study.

With sound motivation and many worked practical examples, the books show in down-to-earth terms how to select and use an appropriate range of statistical techniques in a particular practical field within each title's special topic area.

The books provide statistical support for professionals and research workers across a range of employment fields and research environments. Subject areas covered include medicine and pharmaceutics; industry, finance, and commerce; public services; the earth and environmental sciences, and so on.

The books also provide support to students studying statistical courses applied to the above areas. The demand for graduates to be equipped for the work environment has led to such courses becoming increasingly prevalent at universities and colleges.

It is our aim to present judiciously chosen and well-written workbooks to meet everyday practical needs. Feedback of views from readers will be most valuable to monitor the success of this aim.

A complete list of titles in this series appears at the end of the volume.

To Yea-Jae, Yi,
Tingting, and Shuqin

# PREFACE

With a strong practical emphasis on health science applications, this book describes statistical methods and models for the analysis of data with missing values. We attempt to write so that researchers with experience in applied data analysis, but less technical knowledge than a statistician, should be able to understand and implement most of the methods described. For those with a stronger background in statistics, we provide more technical details as to not detract from the flow of rest of the chapter. We have also tried to choose examples that are relevant to most health science researchers who work in a variety of disciplines.

In all fields of study, missing data are a common problem since, for any data collection process, there are so many things that could go wrong that missing values are all too likely. Thus, when attempts are made to answer the scientific questions of interest, researchers ask the all-too-common question: what do we do with the missing data?

The statistical literature to answer this question is well developed, but overly technical and complicated for researchers who are not experts in statistics and methodology. Therefore, researchers may recognize the existence of missing data, but fail to respond for two reasons: first, they may not understand the consequences of ignoring missing data and how it can impact the validity of their results; second, there is a lack of understanding of the statistical methods for missing data and how to apply them in their own research. Therefore, the purpose of this book is to provide health science researchers with the means of understanding the importance of missing data in their own personal research and the ability to use these methods given the available software.

This book is organized into eight chapters. Chapter 1 introduces concepts on the missing data mechanism and some real-world examples. Chapter 2 gives an overview of methods for dealing with missing data. Chapter 3 describes some design strategies

for minimizing the impact of missing data. Chapters 4 and 5 introduce methods for dealing with missing data problems in cross-sectional and longitudinal studies, respectively. Chapter 6 deals exclusively with missing data problems in survival analysis. Whereas Chapters 3–6 deal with ignorable missing data problems, Chapter 7 presents methods for dealing with nonignorable missing data problems. Finally, Chapter 8 discusses methods for dealing with missing data in causal inferences.

As we worked through examples in the book, we chose to provide software code in the text of the chapters as we want to encourage application of these methods after an understanding of the basic theory. We chose to include R code in the text as many of the methods can be implemented in R; in addition, R is also a publicly available software environment (see www.r-project.org). Since many researchers also use Stata in addition to R, we include code for some selected examples. All the analysis data sets, together with R and Stata codes used in this book, can be downloaded from http://faculty.washington.edu/azhou/.

X.H. ZHOU, C. ZHOU, D. LIU, AND X. DING

*Seattle, Washington*
*March, 2014*

# CONTENTS IN BRIEF

# CONTENTS

# LIST OF FIGURES

# LIST OF TABLES