

SAS for Epidemiologists

Charles DiMaggio

SAS for Epidemiologists

Applications and Methods



Springer

Charles DiMaggio, PhD
Departments of Anesthesiology
and Epidemiology
College of Physicians and Surgeons
Mailman School of Public Health
Columbia University
New York, USA

ISBN 978-1-4614-4853-2 ISBN 978-1-4614-4854-9 (eBook)
DOI 10.1007/978-1-4614-4854-9
Springer New York Heidelberg Dordrecht London

Library of Congress Control Number: 2012947994

© Springer Science+Business Media New York 2013

This work is subject to copyright. All rights are reserved 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. Exempted from this legal reservation are brief excerpts in connection with reviews or scholarly analysis or material supplied specifically for the purpose of being entered and executed on a computer system, for exclusive use by the purchaser of the work. Duplication of this publication or parts thereof is permitted only under the provisions of the Copyright Law of the Publisher's location, in its current version, and permission for use must always be obtained from Springer. Permissions for use may be obtained through RightsLink at the Copyright Clearance Center. Violations are liable to prosecution under the respective Copyright Law.

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.

While the advice and information in this book are believed to be true and accurate at the date of publication, neither the authors nor the editors nor the publisher can accept any legal responsibility for any errors or omissions that may be made. The publisher makes no warranty, express or implied, with respect to the material contained herein.

Printed on acid-free paper

Springer is part of Springer Science+Business Media (www.springer.com)

*For Liz, Katie, and CJ. Thank you for being
there when I come home from epidemiology.*

Foreword

Advances in information technology have transformed the field of epidemiology in profound ways. Data, which had been the most valuable commodity of the empirical sciences, are increasingly available and accessible. The quantity of data rises exponentially on a daily basis. The forms and sources of data are also multiplying. These changes are driving epidemiology into the era of “big data science” and converting more and more epidemiologists into data analysts. As a result, practical skills to manage large and complex data sets and make sound use of them are of increasing importance to future epidemiologists. This book by Dr. Charles DiMaggio is a valuable addition to the toolbox for epidemiology students, public health professionals, and researchers in the health-care industry.

This book distinguishes itself from other applied SAS texts in three notable aspects. First, it is extremely reader-friendly. Dr. DiMaggio has taught the introductory SAS course at Columbia to hundreds of students for over a decade. This book is based primarily on his lectures and is written in a conversational style. The materials are presented in a way that is easy to understand and interesting to learn. In addition to the technical know-how, each chapter contains some “clinical pearls” – insight and wisdom that are unavailable in other applied SAS books. In fact, the reader may feel more like being engaged in a conversation with Dr. DiMaggio than studying a monotonic, mind-numbing how-to manual. Second, this book focuses on the fundamentals. In the first five chapters, Dr. DiMaggio explains in painstaking detail the most basic functions of data input, management, and exploration. These chapters are followed by intermediate statistical analyses of epidemiologic data, both categorical and continuous. Going through these chapters will not make the reader an expert SAS programmer, but it will provide the reader with the necessary skills to perform analytical responsibilities required for a master’s level epidemiologist. Finally, this book is filled with illuminating examples from actual epidemiology research projects. It is written for aspiring epidemiologists by an experienced epidemiologist. These examples are used not only for procedural demonstrations but also for explanations of important epidemiological concepts such as confounding, disease odds ratio, and exposure odds ratio.

As one of the most sophisticated statistical packages, SAS is so kaleidoscopic that first-time users often find it intimidating. Public health students will find this comprehensive text especially helpful for overcoming their initial fears and confusion. It can serve as a textbook for introductory courses on SAS applications as well as a self-study reference for epidemiologists and other health professionals.

New York, NY, USA
Columbia University

Guohua Li
M. Finster

Preface

The path to this book began, as these things so often do, when I was asked to teach a course. In this case, a semester-long class for master's students on epidemiologic analyses using SAS. Over a few years of preparing for and teaching the material I confronted a combination of practical and conceptual considerations that led me to believe that perhaps there was room for another book about SAS.

On a practical level, working with a program like SAS is a skill I consider necessary for all graduating master's epidemiologists. To be honest, the necessity that the program actually be SAS is based on a circular argument. Many employers of epidemiologists use SAS because their current analysts use SAS, and newly minted analysts will compel additional future analysts to use SAS. This reliance on SAS of potential employers of master's-level epidemiology students may change in the future, but my sense is that it will not be anytime soon. While the practical motivation to learn SAS is somewhat self-fulfilling, it does not detract from the capabilities that made SAS an important skill in the first place. And, does it make the choice of SAS any less necessary. As I sit and write this, a quick search on the *New York Times* jobs link returns 15 epidemiology jobs in the New York City area. A search for SAS returns 457 hits. When I do this search on the first day of class, with generally the same results, there is invariably an increased interest among the students in spending a few hours a week learning SAS.

The kinds of SAS-related work that master's-level epidemiologists are called upon to undertake do not exceed some fairly straightforward categorical and continuous data analyses. There was, though, no book that addressed this material in a similarly straightforward fashion. The feedback I've received from the past students is that the procedures covered in this material account for a good majority of their daily activity and that knowing how to do those things helps set the stage for learning more advanced material.

On a conceptual level, the role of statistical software in epidemiologic practice is in a state of flux, and the kinds of data and analyses epidemiologists are being called upon to work with are evolving into what might be called the era of "big data" and the rise of "computational epidemiology." SAS is tailor-made to deal with the kinds of huge data sets that are becoming routine in epidemiology. That there has been

an explosion in the availability of administrative and routinely collected health data, free and open-source data, social media data, and other online data is clear. That the data are amenable to reliable or valid analyses is less clear. The basics, about missing and incorrect values, about confounding, about bias, about study design, are if anything even more important. The data can inform, but we may have to teach them to speak clearly, and in a language that is epidemiologically valid.

Fortunately, SAS is more than up to the task. It has a facility for dealing with extremely large data sets that I have found unsurpassed in other statistical programs. SAS allows epidemiologists to pay special attention to the necessary (though not glamorous) initial steps of reading in, preparing, and cleaning large amounts of data, when early errors or missteps will be amplified throughout the analysis, sometimes in ways that are difficult to trace to their origins. For this reason, fully the first third of this book addresses using SAS to read in and manipulate data to get them into a form that makes epidemiologic sense.

The one aspect of preparing and teaching this material that I did not expect was that it was actually fun. I'm certain this says more about me than it does about the material. But perhaps a kind of geeky enjoyment of some of the practical aspects of epidemiologic methods, like learning how to use SAS, is a sign that you've chosen the right profession. I tried to capture some of what I found interesting and enjoyable by using examples and materials that have practical relevance to epidemiologic practice.

I have come to appreciate that public health practice requires a long-term view, and that you may not always (or even frequently) see the effects of your work. The effects of teaching public health are even farther removed from immediacy. Despite the practical aspects underlying this book, the ultimate motivation is as ephemeral as public health practice itself. In the end, I hope to contribute in some small way to the efforts of someone I haven't met, to improving the health, happiness, and well-being of someone who may not even be born yet.

New York, NY, USA

Charles DiMaggio

Acknowledgments

I was fortunate enough to learn (and hopefully pass on to you) these concepts and procedures from a small but uniformly excellent set of books and the scientists who wrote them [1–12]. I try to cite areas where they were particularly illustrative, but I borrowed from them shamelessly throughout.

I am especially indebted to the SAS Institute for allowing me to draw on their training manuals, course notes, and data sets. If you want a thorough grounding in SAS, I strongly recommend you take one of their outstanding training classes. I also appreciate the willingness of the New York State Department of Health to allow me to use a version of their Statewide Planning and Research Cooperative System (SPARCS) data that figure so prominently in these pages.

Finally, many thanks to my colleagues in the departments of anesthesiology and epidemiology at Columbia University in New York, for their support and feedback, and to my stellar teaching assistants, Joanne Brady, and George Loo, who know SAS so much better than I do, and who helped breathe life into this material.

Contents

1	Introduction	1
1.1	About SAS	1
1.1.1	Alternatives to SAS	2
1.1.2	Why SAS, Then?	3
1.2	About This Book	3
1.2.1	Goals	4
1.2.2	How to Use the Book	4
 Part I Working with Data in SAS		
2	The SAS Environment	9
2.1	The SAS Screen	9
2.2	The Program Editor	9
2.3	SAS Statements	11
2.4	Comments	11
2.5	Quick Demonstration of an SAS Program	12
2.6	Two Types of SAS Programs	13
2.7	Two Kinds of SAS Data	15
2.8	Two Parts to a SAS Data Set	15
2.9	Some Simple SAS Utilities	16
2.10	Getting Help	16
	Problems	17
3	Working with SAS Data	19
3.1	SAS Data Libraries	19
3.2	Two Special SAS Libraries	20
3.3	Three Ways to Browse SAS Data Libraries	21
3.4	Inputting Data into SAS	21
3.5	Reading in Data from the Editor Window	23
3.6	Two Basic INPUT Statements	24
3.6.1	Space-Delimited and Column Input	25

3.7	Reading in Data from External Files	26
3.7.1	The INFILE Statement	26
3.7.2	Formatted INPUT of External Data Sets	28
3.7.3	Informats	29
3.8	Behind the Scenes of a Data Step	30
3.8.1	Deciphering Error Statements.....	31
3.8.2	Error Messages	32
3.8.3	A Few Other Common Errors.....	35
3.9	Notes on Manipulating Data (or How to Tame an Annoying Data Set)	36
3.9.1	Illogically Arrayed Data.....	37
3.10	Data Input Miscellany.....	38
3.11	Importing Excel Spreadsheets	38
	Problems	39
4	Preliminary Procedures	41
4.1	PROC PRINT.....	41
4.2	PROC SORT.....	42
4.3	Enhancing Output: Titles and Footnotes	44
4.4	LABELS	46
4.5	PROC FORMAT and FORMAT	47
4.6	ODS	52
	Problems	54
5	Manipulating Data	57
5.1	The SET Statement.....	57
5.2	Using SET to Define and Create New Variables.....	58
5.2.1	Operations.....	58
5.2.2	Functions	59
5.2.3	Example: Deaths Following the Terrorist Attacks of September 11, 2001	60
5.3	Adding (Concatenating) Data Sets	62
5.3.1	Concatenating the September 11 Data Set	63
5.4	Merging Data Sets Using MERGE – BY	63
5.4.1	SORT Before You MERGE	64
5.4.2	Merging the 9/11 Data Set	65
5.5	Conditional Expressions Using IF-THEN-ELSE	67
5.6	Conditional Expressions Using a Restricting IF Statement	72
5.6.1	Restricting Variables Read into a New Data Set	73
5.7	Conditional Expressions with SAS Dates	73
5.7.1	Using Dates to Subset the 9/11 Data.....	73
	Problems	76

Part II Descriptive and Categorical Analysis

- 6 Descriptive Statistics** 79
 - 6.1 PROC MEANS 79
 - 6.2 PROC FREQ 80
 - 6.3 PROC TABULATE..... 81
 - 6.3.1 Using TABULATE for Surveillance Data 84
 - Problems 86
- 7 Histograms and Plots**..... 91
 - 7.1 Introduction..... 91
 - 7.2 PROC GCHART for Histograms 91
 - 7.3 PROC GPLOT to Plot Continuous Data 93
 - Problems 97
- 8 Categorical Data Analysis I** 99
 - 8.1 Introduction to Categorical Outcomes 99
 - 8.2 Associations 100
 - 8.3 Examining Frequency Tables 101
 - 8.4 Reordering Categorical Variables 103
 - 8.5 Tests of Statistical Significance for Categorical Variables 105
 - 8.5.1 Chi-Square 105
 - 8.5.2 Exact Tests 108
 - 8.5.3 The Mantel–Haenszel Chi-Square 112
 - 8.5.4 The Spearman Correlation Coefficient 113
 - 8.6 Significance vs. Strength..... 114
 - Problems 116
- 9 Categorical Data Analysis II** 119
 - 9.1 Probabilities and Odds 119
 - 9.2 The Odds Ratio 120
 - 9.2.1 Why Epidemiologists Need the Odds Ratio..... 120
 - 9.2.2 The Disease Odds Ratio 121
 - 9.2.3 The Exposure Odds Ratio 123
 - 9.3 Preterm Labor and Birth Weight Example 1 124
 - 9.4 Confounding..... 126
 - 9.4.1 Identifying and Controlling Confounding 126
 - 9.5 Controlling for Confounding 128
 - 9.5.1 Controlling Confounding in Study Design 128
 - 9.5.2 Analytic Approaches to Confounding 128
 - 9.6 Preterm Labor and Birth Weight Example 2 129
 - 9.7 Adjusted Odds Ratios 129
 - 9.7.1 Cochran–Mantel–Haenszel Statistic 131
 - 9.7.2 The Mantel–Haenszel Odds Ratio 131
 - 9.8 Summarizing Exploratory Contingency Table Analyses 135
 - Problems 135

Part III Continuous Data and Regression

10	Cleaning and Assessing Continuous Data using MEANS, UNIVARIATE, and BOXPLOT	139
10.1	PROC MEANS (Redux)	139
10.2	Review of Some Basic Statistics for Continuous Variables	142
10.2.1	Confidence Intervals	147
10.3	PROC UNIVARIATE	148
10.4	PROC BOXPLOT	153
10.5	In Summary	156
	Problems	156
11	ANOVA	159
11.1	Review of ANOVA	159
11.1.1	Assumptions for ANOVA	161
11.2	Testing Assumptions with MEANS, UNIVARIATE, and BOXPLOT	162
11.3	ANOVA with PROC GLM	163
11.3.1	GLM ANOVA Output	166
11.3.2	Multiple Comparisons	167
11.4	Demonstration of One-Way ANOVA	169
11.5	Accounting for More than 1 Categorical Variable: n-Way ANOVA and Interaction Effects	175
11.6	Interaction and Effect Modification: An Epidemiological Perspective	180
11.6.1	The Conundrum of Interaction	180
11.6.2	Components and Causes	182
11.6.3	Usefulness of the Additive Model	183
11.6.4	A Final Thought on Interaction in Epidemiological Studies	184
	Problems	185
12	Correlation	187
12.1	Assessing Correlation	187
12.2	Assessing Correlation Using PROC CORR	188
	Problems	193
13	Linear Regression	197
13.1	Introduction to Regression	197
13.1.1	Variance Perspective of Regression	199
13.2	PROC REG	200
13.2.1	Regression Results	201
13.2.2	Predicted Values	202
13.2.3	Confidence and Prediction Intervals	202
13.3	Demonstration of PROC REG	202

- 13.4 Multiple Regression with PROC REG 206
- 13.5 Interpreting Coefficients 208
 - 13.5.1 Categorical Predictor Variables 208
 - 13.5.2 Demonstration of Multiple Linear Regression 211
- Problems 212
- 14 Regression Diagnostics** 213
 - 14.1 Introduction 213
 - 14.2 Residuals Redux 214
 - 14.2.1 Residual Plots 216
 - 14.3 Outliers (Influential Observations) 217
 - 14.4 Collinearity 217
 - 14.5 Demonstration: Residual Diagnostics for the Fitness Data 219
 - 14.6 A Word About Model Selection 225
 - 14.6.1 SAS Model Selection Tools 225
 - 14.6.2 Problems with Automated Selection Procedures 227
 - 14.6.3 Some Advice on Model Selection 228
 - Problems 229
- References** 231
- Solutions** 233
- Index** 253

