

# Quantitative Methods in Population Health

Extensions of Ordinary Regression

Mari Palta



A JOHN WILEY & SONS, INC., PUBLICATION



# Quantitative Methods in Population Health

WILEY SERIES IN PROBABILITY AND STATISTICS

Established by WALTER A. SHEWHART and SAMUEL S. WILKS

Editors: *David J. Balding, Peter Bloomfield, Noel A. C. Cressie,  
Nicholas I. Fisher, Iain M. Johnstone, J. B. Kadane, Louise M. Ryan,  
David W. Scott, Adrian F. M. Smith, Jozef L. Teugels;*  
Editors Emeriti: *Vic Barnett, J. Stuart Hunter, David G. Kendall*

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

# Quantitative Methods in Population Health

Extensions of Ordinary Regression

Mari Palta



A JOHN WILEY & SONS, INC., PUBLICATION

This book is printed on acid-free paper.②

Copyright © 2003 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](http://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, e-mail: [permreq@wiley.com](mailto:permreq@wiley.com).

**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 please contact our Customer Care Department within the U.S. at 877-762-2974, outside the U.S. 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, however, may not be available in electronic format.

For ordering and customer service, call 1-800-CALL-WILEY.

***Library of Congress Cataloging-in-Publication Data:***

Palta, Mari, 1948–

Quantitative methods in population health: extensions of ordinary regression / Mari Palta.

p. cm.

Includes bibliographical references and index.

ISBN 0-471-45505-9 (cloth)

1. Medical statistics. 2. Regression analysis. 3. Population–Health aspects–Statistical methods. 4. Health surveys–Statistical methods. I. Title.

RA409.P34 2003

614.4'2'0727—dc21

2003050087

Printed in the United States of America.

10 9 8 7 6 5 4 3 2 1

# List of Figures

<i>1.1</i>	<i>The number of participants with GHb</i>	<i>xxiii</i>
<i>1.1</i>	<i>Wisconsin Diabetes Registry Study GHb for the first five subjects</i>	<i>8</i>
<i>7.1</i>	<i>Missing data in describing the relationship between health and exercise</i>	<i>108</i>





# List of Tables

6.1	<i>Regression coefficients (se) modeling GHb to ages up to 15 years</i>	65
6.2	<i>Regression coefficients (se) modeling SBP in Wisconsin Sleep Cohort</i>	66
6.3	<i>Regression coefficients (se) modeling transformed GHb on age</i>	67
6.4	<i>Regression coefficients (se) modeling GHb to ages up to 15 Years by weighted regression</i>	72
7.1	<i>Unweighted and weighted <math>\hat{\beta}</math> (empirical se) of Wisconsin Sleep Cohort SBP</i>	102
7.2	<i>Comparison of birth characteristics between participants and nonparticipants in functional assessment at age 5 years</i>	104
7.3	<i>Comparison of birth characteristics between all survivors and participants in functional assessment at age 5 years</i>	105
7.4	<i>Unweighted and weighted regression coefficients (empirical se) for predicting social function at age 5 years for VLBW children</i>	105
8.1	<i>Regression results for all visits <math>\hat{\beta}</math> (se) and empirical se</i>	124
8.2	<i><math>\hat{\beta}</math> (model-based se), (empirical se) based on different correlation structures</i>	133
8.3	<i><math>\hat{\beta}</math> (model-based se), (empirical se) based on autoregressive AR(1) variance structure</i>	134
9.1	<i><math>\hat{\beta}</math> (empirical se) fitting an overall model by compound symmetry, and fitting separate between- and within-individual effects for age</i>	153
10.1	<i><math>\hat{\beta}</math> (se) and variance components (se) random intercept model</i>	160
10.2	<i><math>\hat{\beta}</math> (se) from a model with random slope and from an AR(1) model</i>	166
12.1	<i>Characteristics of three distributions in the exponential family</i>	192
14.1	<i>Standard errors for Wisconsin cancer mortality in 1996 model, taking overdispersion into account</i>	246

14.2	<i>Standard errors for infant hospitalization model, taking overdispersion into account</i>	246
15.1	<i>Comparing results from PROC MIXED and PROC GENMOD for SBP</i>	268
15.2	<i><math>\hat{\beta}</math>(model se) for independence and <math>\hat{\beta}</math> (empirical se) for compound symmetry</i>	270
15.3	<i>Adjusting for number of follow-up visits and for survey weights</i>	270
15.4	<i>Risk, odds ratios, and 95% confidence intervals from final model</i>	271
15.5	<i>Regression coefficients (empirical se) for hospitalization of VLBW children</i>	272
15.6	<i>Rate and rate ratios (empirical se) for hospitalization of VLBW children</i>	273
15.7	<i>Odds ratios [95% CI] from GEE and from conditional logistic regression</i>	274

# Contents

<b>Preface</b>	<b>xv</b>
<b>Acknowledgments</b>	<b>xvii</b>
<b>Acronyms</b>	<b>xix</b>
<b>Introduction</b>	<b>xxi</b>
I.1    Newborn Lung Project	xxi
I.2    Wisconsin Diabetes Registry	xxii
I.3    Wisconsin Sleep Cohort Study	xxiii
<b>Suggested Reading</b>	<b>xxv</b>
<b>1    Review of Ordinary Linear Regression and Its Assumptions</b>	<b>1</b>
1.1    The Ordinary Linear Regression Equation and Its Assumptions	1
1.1.1    Straight-Line Relationship	2
1.1.2    Equal Variance Assumption	5
1.1.3    Normality Assumption	6
1.1.4    Independence Assumption	7
1.2    A Note on How the Least-Squares Estimators are Obtained	8
Output Packet I: Examples of Ordinary Regression Analyses	9
<b>2    The Maximum Likelihood Approach to Ordinary Regression</b>	<b>21</b>
2.1    Maximum Likelihood Estimation	21
2.2    Example	24
2.3    Properties of Maximum Likelihood Estimators	25
2.4    How to Obtain a Residual Plot with PROC MIXED	26
Output Packet II: Using PROC MIXED and Comparisons to PROC REG	26
	<b>ix</b>

<b>3</b>	<b>Reformulating Ordinary Regression Analysis in Matrix Notation</b>	<b>30</b>
3.1	Writing the Ordinary Regression Equation in Matrix Notation	31
3.1.1	Example	32
3.2	Obtaining the Least-Squares Estimator $\hat{\beta}$ in Matrix Notation	33
3.2.1	Example: Matrices in Regression Analysis	34
3.3	List of Matrix Operations to Know	36
<b>4</b>	<b>Variance Matrices and Linear Transformations</b>	<b>38</b>
4.1	Variance and Correlation Matrices	38
4.1.1	Example	40
4.2	How to Obtain the Variance of a Linear Transformation	40
4.2.1	Two Variables	40
4.2.2	Many Variables	42
<b>5</b>	<b>Variance Matrices of Estimators of Regression Coefficients</b>	<b>51</b>
5.1	Usual Standard Error of Least-Squares Estimator of Regression Slope in Nonmatrix Formulation	51
5.2	Standard Errors of Least-Squares Regression Estimators in Matrix Notation	52
5.2.1	Example	53
5.3	The Large Sample Variance Matrix of Maximum Likelihood Estimators	54
5.4	Tests and Confidence Intervals	56
5.4.1	Example-Comparing PROC REG and PROC MIXED	57
<b>6</b>	<b>Dealing with Unequal Variance Around the Regression Line</b>	<b>62</b>
6.1	Ordinary Least Squares with Unequal Variance	62
6.1.1	Examples	64
6.2	Analysis Taking Unequal Variance into Account	66
6.2.1	The Functional Transformation Approach	66
6.2.2	The Linear Transformation Approach	68
6.2.3	Standard Errors of Weighted Regression Estimators	73
	Output Packet III: Applying the Empirical Option to Adjust Standard Errors	75
	Output Packet IV: Analyses with Transformation of the Outcome Variable to Equalize Residual Variance	83
	Output Packet V: Weighted Regression Analyses of GHb Data on Age	93

<b>7 Application of Weighting with Probability Sampling and Nonresponse</b>	<b>97</b>
7.1 Sample Surveys with Unequal Probability Sampling	98
7.1.1 Example	101
7.2 Examining the Impact of Nonresponse	102
7.2.1 Example (of Reweighting as Well as Some SAS Manipulations)	104
7.2.2 A Few Comments on Weighting by a Variable Versus Including it in the Regression Model	107
Output Packet VI: Survey and Missing Data Weights	109
<b>8 Principles in Dealing with Correlated Data</b>	<b>119</b>
8.1 Analysis of Correlated Data by Ordinary Unweighted Least-Squares Estimation	120
8.1.1 Example	121
8.1.2 Deriving the Variance Estimator	122
8.1.3 Example	124
8.2 Specifying Correlation and Variance Matrices	124
8.3 The Least-Squares Equation Incorporating Correlation	126
8.3.1 Another Application of the Spectral Theorem	127
8.4 Applying the Spectral Theorem to the Regression Analysis of Correlated Data	128
8.5 Analysis of Correlated Data by Maximum Likelihood	129
8.5.1 Non equal Variance	130
8.5.2 Correlated Errors	131
8.5.3 Example	132
Output Packet VII: Analysis of Longitudinal Data in Wisconsin Sleep Cohort	135
<b>9 A Further Study of How the Transformation Works with Correlated Data</b>	<b>145</b>
9.1 Why Would $\beta_W$ and $\beta_B$ Differ?	147
9.2 How the Between- and Within-Individual Estimators are Combined	149
9.3 How to Proceed in Practice	151
9.3.1 Example	152
Output Packet VIII: Investigating and Fitting Within- and Between-Individual Effects	154
<b>10 Random Effects</b>	<b>156</b>
10.1 Random Intercept	156
10.1.1 Example	159

10.1.2	Example	161
10.2	Random Slopes	161
10.2.1	Example	165
10.3	Obtaining “The Best” Estimates of Individual Intercepts and Slopes	167
10.3.1	Example	167
	Output Packet IX: Fitting Random Effects Models	169
<b>11</b>	<b>The Normal Distribution and Likelihood Revisited</b>	<b>181</b>
11.1	PROC GENMOD	182
11.1.1	Example	183
	Output Packet X: Introducing PROC GENMOD	184
<b>12</b>	<b>The Generalization to Non-normal Distributions</b>	<b>190</b>
12.1	The Exponential Family	190
12.1.1	The Binomial Distribution	192
12.1.2	The Poisson Distribution	193
12.1.3	Example	194
12.2	Score Equations for the Exponential Family and the Canonical Link	194
12.3	Other Link Functions	196
12.3.1	Example	197
<b>13</b>	<b>Modeling Binomial and Binary Outcomes</b>	<b>199</b>
13.1	A Brief Review of Logistic Regression	199
13.1.1	Example: Review of the Output from PROC LOGIST	200
13.2	Analysis of Binomial Data in the Generalized Linear Models Framework	202
13.2.1	Example of Logistic Regression with Binary Outcome	206
13.2.2	Example with Binomial Outcome	207
13.2.3	Some More Examples of Goodness-of-Fit Tests	209
13.3	Other Links for Binary and Binomial Data	209
13.3.1	Example	211
	Output Packet XI: Logistic Regression Analysis with PROC LOGIST and PROC GENMOD	212
	Output Packet XII: Analysis of Grouped Binomial Data	221
	Output Packet XIII: Some Goodness-of-Fit Tests for Binomial Outcome	223
	Output Packet XIV: Three Link Functions for Binary Outcome	229
	Output Packet XV: Poisson Regression	247
	Output Packet XVI: Dealing with Overdispersion in Rates	254

<b>14 Modeling Poisson Outcomes—The Analysis of Rates</b>	<b>236</b>
14.1 Review of Rates	236
14.1.1 Relationship Between Rate and Risk	238
14.2 Regression Analysis	239
14.3 Example with Cancer Mortality Rates	241
14.3.1 Example with Hospitalization of Infants	242
14.4 Overdispersion	243
14.4.1 Fitting a Dispersion Parameter	244
14.4.2 Fitting a Different Distribution	245
14.4.3 Using Robust Standard Errors	245
14.4.4 Applying Adjustments for Over Dispersion to the Examples	246
Output Packet XV: Poisson Regression	247
<b>15 Modeling Correlated Outcomes with Generalized     Estimating Equations</b>	<b>263</b>
15.1 A Brief Review and Reformulation of the Normal Distribution, Least Squares and Likelihood	263
15.2 Further Developments for the Exponential Family	264
15.3 How are the Generalized Estimating Equations Justified?	266
15.3.1 Analysis of Longitudinal Systolic Blood Pressure by PROC MIXED and GENMOD	267
15.3.2 Analysis of Longitudinal Hypertension Data by PROC GENMOD	269
15.3.3 Analysis of Hospitalizations Among VLBW Children Up to Age 5	271
15.4 Another Way to Deal with Correlated Binary Data	273
Output Packet XVII: Mixed Versus GENMOD for Longitudinal SBP and Hypertension Data	274
Output Packet XVIII: Longitudinal Analysis of Rates	285
Output Packet XIX: Conditional Logistic Regression of Hypertension Data	288
<b>References</b>	<b>290</b>
<b>Appendix: Matrix Operations</b>	<b>295</b>
A.1 Adding Matrices	296
A.2 Multiplying Matrices by a Number	297
A.3 Multiplying Matrices by Each Other	297
A.4 The Inverse of a Matrix	299
<b>Index</b>	<b>303</b>





# Preface

This text arose from my many years working with several long-term population-based observational studies. As I was asked to put together a third-semester statistics course for our new Ph.D. program in Population Health, I decided to assemble the information I had seen investigators and students need most often, and I also decided to answer as many questions as possible out of those I had typically been asked. The resulting mix of topics is guided by this experience. I have attempted to pull in and deal with the aberrations of observational data such as confounding and selection bias. Some traditional topics regarding small sample inference, analysis of variance, and experimental design are deemphasized, as I have found that they confuse rather than help population health researchers. I am using data sets from my own research and collaborations as examples to ensure that subject matter interpretations are meaningful, and that the reader becomes familiar with the “non-textbook” appearance of real data.

While keeping the material immediately applicable by providing detailed instructions for how to run and interpret procedures in SAS, I find it irresponsible to do so without creating some “common sense” about the methods and their assumptions. The beginning chapters lay the mathematical groundwork necessary for topics in later chapters. Whenever possible, I have made a point of inserting practical issues that are answered by specific mathematical derivations.

In addition, each topic starts with an explanation of the theoretical background necessary to allow reasonable judgment as to when the technique is applicable and to facilitate future learning of related methods and software. In the process the reader is exposed to some of the underpinnings of statistics that are often omitted from applied texts and courses. While the text is anchored in the terminology of the biostatistical tradition, I point to some important connections to techniques and terminology used in econometrics and psychometrics. Because of the historic emphasis of biostatistics on experiments and randomization, I have often found that econometric approaches provide further insight in the observational framework.

For progress in addressing current population health issues, it is necessary for researchers in epidemiology and health services to understand and apply regression analysis with weights to deal with unequal variance and correlated and longitudinal outcomes by mixed effects, generalized linear models, and generalized estimating equations. In addition, many data sets in these areas include survey weights. Increasingly, investigators are also called upon to examine the possible impact on their results of missing observations. I suggest straightforward methodology that can be implemented with standard software. The material is presented on a level that will make it accessible to epidemiologists and health services researchers, as well as to applied statisticians. This corresponds roughly to a third-semester applied statistics sequence for statistics non-majors. It is assumed that the reader is already well acquainted with ordinary and logistic regression analysis and has at least rudimentary knowledge of the SAS package.

The explanations are designed to assume as little background in mathematics and statistical theory as possible, except that some knowledge of calculus is necessary for certain parts, such as in understanding maximum likelihood and generalized linear models. The reader may wish to review the rules and uses of derivatives, which are not covered here. On the other hand, all relevant aspects of linear algebra and statistical theory are explained within the text. Important formulas are derived, but with an eye to avoiding excessive algebra.

SAS commands are provided for applying the methods. The SAS procedures emphasized are PROC REG, PROC MIXED, and PROC GENMOD, with occasional references to others. Useful data manipulation commands are introduced as needed to illustrate the techniques in the specific data sets, and the SAS ODS system is briefly introduced to accomplish viewing random effects from mixed models. However, basic commands used to read in data sets and annotate them are considered well known and are not always provided in the text.

Mari Palta  
*Madison, Wisconsin*

# Acknowledgments

I have been fortunate to collaborate with wonderful investigators in epidemiology and health services research. They have continually inspired me to look into new statistical issues and have provided me with the insights needed to link the statistics to the subject matter. My longest and closest collaborations have been with Mona Sadek-Badawi in research on the outcomes of very-low-birth-weight neonates, with Kit Allen in diabetes research, and with Terry Young in research in sleep disorders. We have together tackled many challenges that arise in long-term subject follow-up, such as the difficulty in locating subjects and the effects of some participants skipping samples or questionnaires. The research has also interested me, among many other things, in latent variables and in examining between and within individual effects. Many other individuals working on these projects, including Tammy LeCaire and Laurel Finn, have helped me obtain the data for the examples.

Many of the specifics addressed in this book arose from the research of Paul Peppard on the association of sleep apnea with blood pressure. These data form a core example that is traced through the various analytic techniques. I appreciate the many discussions and penetrating questions raised by Paul. Among health services researchers, Maureen Smith has perked my interest in econometric and other social science techniques and has lent me many books. Dennis Fryback shared his data and questions on health care cost analyses.

The book could not have been pulled together without the help of Lin Wang, a Ph.D. student in Statistics who critically reviewed the text, served as a sounding board for methodology, and helped with the formatting of the text. Other former statistics Ph.D. students whose research contributed to the methodology are T.-Y. Yao, Chin-Yu Lin, Wei-Hsiung Chao, Soomin Park, Lei Shen and Liang Li. Many students who have taken the class Quantitative Methods in Population Health, for which this text was initially drafted, have challenged me to explain statistics more clearly and have forced me to consider how different techniques relate to each other and why we do statistics the way we do.

Last but not least, I owe a lot to all the people who have supported me personally and who have pitched in when and wherever needed. My family has been wonderfully accommodating to my obsession with this book. I could not have completed this task without their love and support. In addition, everyone working on the Newborn Lung Project, including Aggie Albanese, Kathleen Madden, and Hana Said and the staff of the Diabetes Registry, make my immediate work environment fun and rewarding.

New methodological developments underlying certain parts of the book were supported by grant CA53786 from the National Cancer Institute and by grant P01 HL42242 from the National Heart Lung and Blood Institute.

M. P.

# Acronyms

ANOVA	Analysis of variance
BLUE	Best linear unbiased estimator
GHb	Glycosylated hemoglobin
ML	Maximum likelihood
NICU	Neonatal intensive care unit
NHANES	National health and nutrition examination survey
REML	Restricted maximum likelihood
SAS	Statistical Analysis System®
SBP	Systolic blood pressure
se	standard error
VLBW	Very low birth weight



# Introduction

## **Some Data Sets Used as Examples in This Text**

In this book, we focus on extending ordinary regression analysis by considering situations where some of the usual assumptions are violated. As we discuss in more detail in future chapters, violations of assumptions are common in population health data sets. For example, when we want to model presence versus absence of a disease the outcome variable is binary. Because ordinary regression is designed for normally distributed outcomes, the presence of binary outcomes leads to the extension of ordinary regression to logistic regression. We discuss this extension and others that apply to non-normally distributed outcomes in later chapters. With normally distributed outcomes, we encounter situations with violation of equal variance and independence assumptions. For example, subjects may be followed longitudinally, which leads to correlated residuals. The fitting of models and inference in such cases is the topic of the earlier chapters. Throughout, we illustrate with the use of data sets that have accrued from population health research. The purpose of this introduction is to briefly describe these data sets. All analyses were run in SAS 8.2 [1] for Solaris. For graphics we occasionally used SAS 8.2 for Windows.

## **1.1 NEWBORN LUNG PROJECT**

The Newborn Lung Project enrolled a cohort that included all very-low-birth-weight admissions to six neonatal intensive care units in Wisconsin and Iowa during 8/1/88–6/30/91. There were 1040 admissions during this time period, and some baseline data were collected on all of them. Neonatal nurses collected medical record information on factors such as birth weight, supplemental oxygen use the first 24 hours, and hospital of birth without identifiers. Parents were approached as soon as possible for informed consent for interview and medical record abstracting. Due to human subjects concerns, parents were not approached after the neonate had

died or if the neonate was in critical condition. A total of 810 infants survived the hospitalization, and the parents of 633 provided informed consent for abstracting. Recontract addresses were available only for the subgroup with informed consent. By age 5, six additional children had died. Among the 804 survivors, 438 were located, and a follow-up interview including health information and a functional assessment of 422 was performed. The parents of 345 children also gave informed consent for complete abstracting of medical records.

The original purpose of the study was to establish severity scores for neonatal lung disease and to find risk factors associated with it. Later, we described functional and respiratory outcomes at ages 5 and 8 years and their predictors. The record abstracting led to longitudinal data on number of hospitalizations and clinic visits during every year of life.

Examples used in the text arise from data collected during the initial hospitalization and from the follow-up. For example, we analyze functional outcome at age 5 and hospitalizations during the first five years of life as outcomes in regression analysis. We briefly illustrate principal component analysis of socioeconomic indicators collected by this study. We will also use this data set to show how to use the available data on those not followed to examine selection bias.

Some references that present data included in this text and that provide further background on the Newborn Lung Project are listed at the end of this Introduction.

## **I.2 WISCONSIN DIABETES REGISTRY**

The Wisconsin Diabetes Registry targeted all individuals <30 years of age diagnosed with Type I diabetes in a 28-county area in southern Wisconsin 5/1/87–4/30/92. The primary mode of recruitment was by physician, diabetes educator, and self-referral. Also, all hospitals and most multipractice clinics were telephoned every 3 months to ascertain any unreferred cases. A total of 733 cases were found. Out of these, 597 gave informed consent for participation. Participants underwent a baseline interview, were requested to submit blood samples every 4 months, were sent a questionnaire inquiring about hospitalizations and other events every 6 months, and underwent physical examinations at 4, 7, and 9 years of duration. The blood samples were used to determine glycosylated hemoglobin (GHb), an important indicator of glycemic control. The purpose of the study is to map the acute and chronic outcomes of Type 1 diabetes from diagnosis and to consider risk factors such as GHb from the earliest stages of the disease onwards.

Examples used in this text arise from the longitudinal glycosylated hemoglobin measurements performed on the blood samples. Figure I.1 obtained by PROC CHART shows the number of participants with GHb data for 1 year, 2 years, and so on, up to 14 years of follow-up, at the time the data sets for this text were compiled. In most analysis here, we will for simplicity average all the GHb measures in a given year. The commands used to produce the number of years of GHb measurements for each individual were:



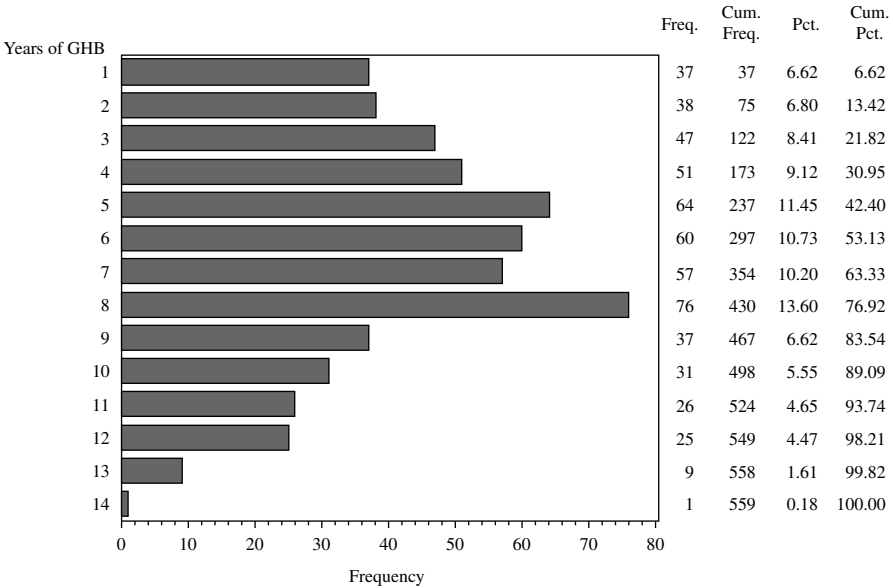


Fig. I.1 The number of participants with GHB

```
PROC SORT; BY ID;
PROC MEANS NOPRINT; BY ID; VAR GHB;
OUTPUT OUT=MM N=NG;
DATA B; SET MM;
LABEL NG='YEARS OF GHB';
PROC CHART; HBAR NG;
```

Some references pertaining to GHb measurements in the Wisconsin Diabetes Registry Study are listed at the end of this Introduction.

I.3 WISCONSIN SLEEP COHORT STUDY

A survey questionnaire inquiring about sleep and sleep-related problems was sent to 6900 employees age 30–60 at four State of Wisconsin agencies. Completed surveys were received from 4927 respondents. A stratified random sample of the respondents was invited to spend the night in a completely equipped clinical sleep laboratory for overnight polysomnography and other tests. A total of 1370 individuals participated. Sleep studies were performed over an extended time period, resulting in some individuals being age 65 and older at the first visit. Blood pressure measurements were taken in the laboratory, and height and weight were measured. Subjects are reinvited for sleep studies every four years. The goals of the project are to identify risk factors and outcomes associated with sleep disorders.

Data used in this text are the longitudinal measures of systolic blood pressure and hypertension as associated with age, gender, and body mass index. Sample sizes in specific analyses vary slightly due to measurements sometimes being missing. For example, out of the 1370 total individuals, 5 had their first blood pressure measurement at the second visit. We also use data from a subproject on the cost of medical care for the cohort. A list of references from the study that involve these variables is provided below.

## Suggested Reading

1. Palta M, Gabbert D, Weinstein MR, and Peters ME, Multivariate assessment of traditional risk factors for chronic lung disease in very low birth weight neonates, *Journal of Pediatrics*, **119**:285–292 (1991).
2. Palta M, Weinstein MR, McGuinness G, Gabbert D, Brady W, and Peters ME, A population study: Mortality and morbidity after availability of surfactant therapy, *Archives of Pediatrics and Adolescent Medicine*, **148**:1295–1301 (1994).
3. Palta M, Sadek M, Barnet JH, Evans M, Weinstein MR, McGuinness G, Peters ME, Gabbert D, Fryback D, and Farrell P, Evaluation of criteria for chronic lung disease in very low birth weight infants, *Journal of Pediatrics*, **132**:57–63 (1998).
4. Palta M, Sadek M, Evans M, Weinstein MR, and McGuinness G, Functional assessment of a multicenter VLBW cohort at age 5 years”, *Archives of Pediatrics and Adolescent Medicine*, **154**:23–30 (2000).

### Wisconsin Diabetes Registry

5. Allen C, Zaccaro D, Palta M, et al., Glycemic control in the first two years of insulin-dependent diabetes mellitus, *Diabetes Care*, **15**:980–987 (1992).
6. Palta M, Shen G, Allen C, Klein R, and D’Alessio D, Longitudinal glycosylated hemoglobin patterns from diagnosis in a population based cohort with IDDM, *American Journal of Epidemiology*, **114**:954–961 (1996).
7. Allen C, LeCaire T, Palta M, Daniels K, Meredith M, and D’Alessio D, Risk factors for frequent and severe hypoglycemia in Type I diabetes, *Diabetes Care*, **24**:1878–1881 (2001).

### Wisconsin Sleep Cohort Study

8. Young T, Palta M, Dempsey J, Skatrud J, Weber S, and Badr S, Occurrence of sleep disordered breathing among middle-aged adults, *New England Journal of Medicine*, **328**:69–77 (1993).

9. Hla KM, Young TB, Bidwell T, Palta M, Skatrud J, and Dempsey J, Sleep apnea and hypertension—A population-based study, *Annals Internal Medicine*, **120**:382–388 (1994).
10. Young T, Finn L, Hla M, Morgan BJ, Palta M, Snoring as part of a dose–response relationship between sleep-disordered breathing and blood pressure, *Sleep*, **19** (supplement) (1996).
11. Peppard PE, Young T, Palta M, and Skatrud J, Association between sleep disordered breathing and hypertension, *New England Journal of Medicine*, **342**:1378–1384 (2000).
12. Young T, Peppard P, Palta M, Hla KM, Finn L, Morgan B, and Skatrud J, Population-based study of sleep disordered breathing as a risk factor for hypertension, *Archives of Internal Medicine*, **157**:1746–1752 (1997).

## CHAPTER ONE

# Review of Ordinary Linear Regression and Its Assumptions

### 1.1 THE ORDINARY LINEAR REGRESSION EQUATION AND ITS ASSUMPTIONS

A linear regression equation can be alternatively specified as

$$\begin{aligned}y_i &= \beta_0 + \beta_1 x_i + \epsilon_i \quad \text{or} \\ \mu_{y|x} &= \beta_0 + \beta_1 x \quad \text{or} \\ E(y|x) &= \beta_0 + \beta_1 x\end{aligned}\tag{1.1}$$

to describe the quantitative relationship between a single predictor  $x$  and an outcome  $y$ . In the population health research projects described in the Introduction,  $y$  may be a measured GHb or the score of a very-low-birth-weight (VLBW) child on a test, or the systolic blood pressure (SBP) at a visit to the sleep clinic. In the first equation  $\epsilon_i$  is a random regression error describing the deviation of a given value  $y_i$  from its mean. It can be viewed as capturing unmeasured influence on the outcome. In order to make both the first and the second equations of (1.1) correct, it is assumed that  $E(\epsilon_i|x_i) = 0$ . In other words, if the second equation is to describe the relationship of the mean  $y$  to  $x$  correctly, the random errors in the first equation must average to 0 for all  $x$ . This also implies that  $\epsilon_i$  does not depend on  $x_i$ . The last two equations are just saying the same thing in different notation because the “expected value”  $E(\cdot)$  of a variable is by definition the mean of that variable.

We assume that the reader is familiar with the “conditional on” notation implied by the “[|]”. Conditioning on a variable means that the variable is (at that moment) considered a constant, so the parameters of the distribution of  $y$  may depend on  $x$ . In other words, when conditioning systolic blood pressure on a given age  $x$ , we are interested in the parameters of the distribution of blood pressure at that age. Estimation of the parameters of equations (1.1) usually proceeds by the method

of least squares. In dealing with the regression equation, forming estimators, and drawing inference, we commonly make a number of assumptions:

### 1.1.1 Straight-Line Relationship

Equation (1.1) implies that  $x$  and the mean of  $y$  are related in a straight-line fashion. This assumption can be alternatively stated as a constant difference in mean  $y$  between every pair of  $x$ 's that are separated by the same number of steps. For example, if  $y$  is systolic blood pressure from visit 1 in the Sleep Cohort Study and  $x$  is age, linearity implies that the difference in mean blood pressure between a 50-year-old and a 40-year-old is the same as that between a 40-year-old and a 30-year-old. Regardless of the level of  $x$ ,  $\mu_{y|x+1} - \mu_{y|x} = \beta_1$ , so that the regression coefficient is the difference in mean with one step increase in  $x$ . Again, if  $y$  is systolic blood pressure and age  $x$  is recorded in years,  $\beta_1$  is the increase in mean blood pressure every year. The linearity assumption is an inherent structural assumption, the validity of which is driven by the biological, sociological, and so on, mechanisms that relate  $y$  to  $x$ . When the linearity assumption holds, we are ahead statistically, because we need to estimate only two parameters  $\beta_0$  and  $\beta_1$  instead of a separate  $\mu_{y|x}$  for every  $x$ .

Only in the situation that  $x$  is binary (e.g., designating two treatment groups) is the linearity assumption moot, or automatically satisfied. If  $y$  is systolic blood pressure and  $x$  is a 0–1 indicator of gender where 1 indicates male, then  $\beta_1$  is the difference in mean blood pressure between males and females, and  $\beta_0$  is the mean for females. In this situation,  $\mu_{y|x}$  is simply a notation for representing the means of two groups (females and one-step difference involved). Since no assumptions are made on the mean structure, equations (1.1) estimate two parameters either way.

In other situations, the original  $x$  may just serve as a label for different groups, such as ethnic categories or treatments. The linearity assumption then makes little sense. However, we can expand (1.1) through the device of binary indicator variables, which bypass the linearity assumption, but again do not save us parameters as compared to estimating  $\mu_{y|x}$  separately for each group. In the Wisconsin Sleep Cohort Study, we may wish to compare mean blood pressure between the four state agencies surveyed, by using three indicator variables. In SAS, indicator variables are created in many procedures by the CLASS statement [1].

In the simple cases presented in this chapter, we emphasize linearity of  $\mu_{y|x}$  versus a single predictor. We can easily generalize equation (1.1) to more complicated cases by transforming  $y$  or  $x$  or by adding squared, cubic, and so on, terms in  $x$ . Note, however, that even when  $x$  or  $y$  is transformed or when polynomial terms are added, ordinary regression remains a linear expression of the regression parameters. This simplifies estimation. In Chapter 12, we will consider some situations when the regression equation for the mean is not linear in the parameters.

#### 1.1.1.1 Example

OUTPUT PACKET I shows regression equations, plots of residuals versus predicted values, and mean plots for some variables from the data sets of interest. Later, we