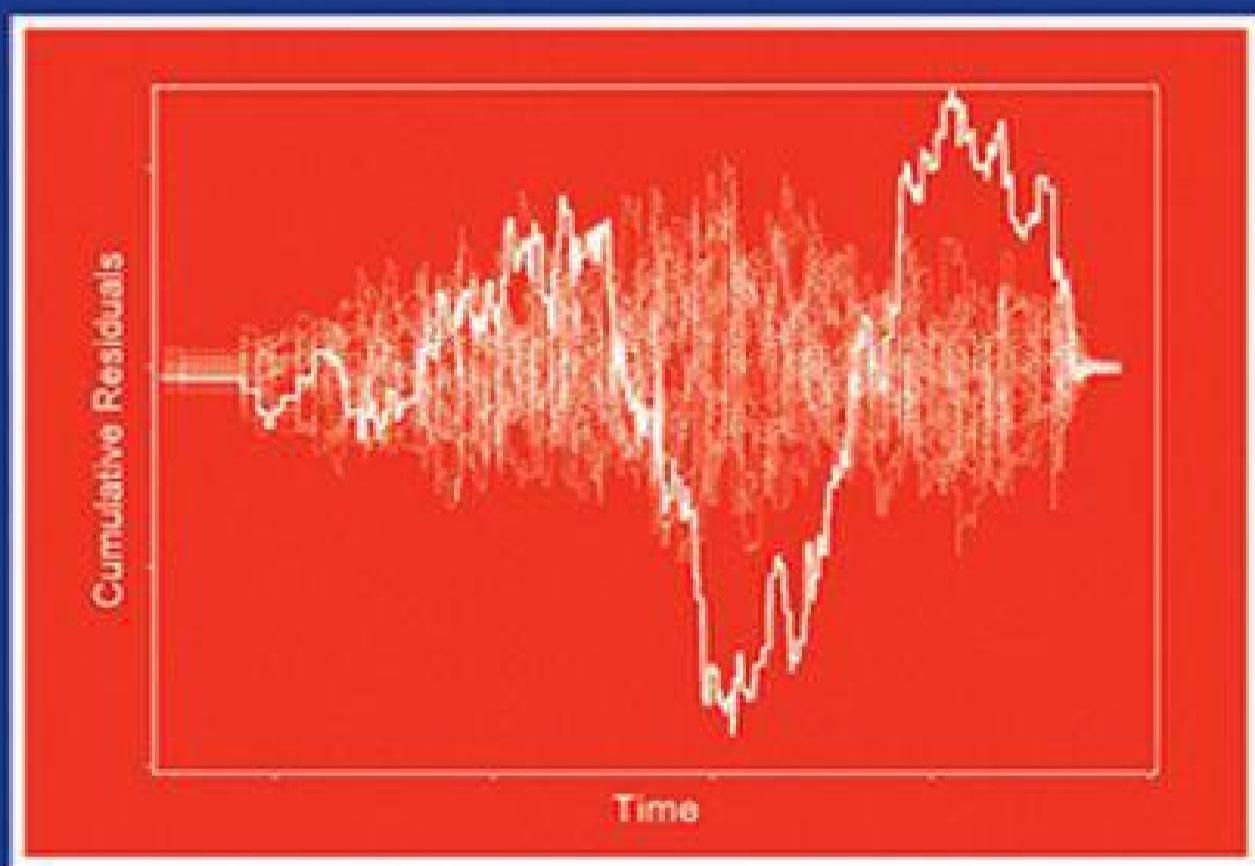


WILEY SERIES IN PROBABILITY AND STATISTICS

# Applied Longitudinal Analysis

## Second Edition



**Garrett M. Fitzmaurice**  
**Nan M. Laird**  
**James H. Ware**

# Contents

[Half Title page](#)

[Title page](#)

[Copyright page](#)

[Dedication](#)

[Preface](#)

[Preface to First Edition](#)

[Acknowledgments](#)

## [Part I: Introduction to Longitudinal and Clustered Data](#)

### [Chapter 1: Longitudinal and Clustered Data](#)

- [1.1 Introduction](#)
- [1.2 Longitudinal and Clustered Data](#)
- [1.3 Examples](#)
- [1.4 Regression Models for Correlated Responses](#)
- [1.5 Organization of the Book](#)
- [1.6 Further Reading](#)

### [Chapter 2: Longitudinal Data: Basic Concepts](#)

- [2.1 Introduction](#)
- [2.2 Objectives of Longitudinal Analysis](#)
- [2.3 Defining Features of Longitudinal Data](#)
- [2.4 Example: Treatment of Lead-Exposed Children Trial](#)
- [2.5 Sources of Correlation in Longitudinal Data](#)
- [2.6 Further Reading](#)

## [Part II: Linear Models for Longitudinal Continuous Data](#)

### [Chapter 3: Overview of Linear Models for Longitudinal Data](#)

- [3.1 Introduction](#)
- [3.2 Notation and Distributional Assumptions](#)
- [3.3 Simple Descriptive Methods of Analysis](#)

- [3.4 Modeling the Mean](#)
- [3.5 Modeling the Covariance](#)
- [3.6 Historical Approaches](#)
- [3.7 Further Reading](#)

## **Chapter 4: Estimation and Statistical Inference**

- [4.1 Introduction](#)
- [4.2 Estimation: Maximum Likelihood](#)
- [4.3 Missing Data Issues](#)
- [4.4 Statistical Inference](#)
- [4.5 Restricted Maximum Likelihood \(REML\) Estimation](#)
- [4.6 Further Reading](#)

## **Chapter 5: Modeling the Mean: Analyzing Response Profiles**

- [5.1 Introduction](#)
- [5.2 Hypotheses Concerning Response Profiles](#)
- [5.3 General Linear Model Formulation](#)
- [5.4 Case Study](#)
- [5.5 One-Degree-of-Freedom Tests for Group by Time Interaction](#)
- [5.6 Adjustment for Baseline Response](#)
- [5.7 Alternative Methods of Adjusting for Baseline Response](#)
- [5.8 Strengths and Weaknesses of Analyzing Response Profiles](#)
- [5.9 Computing: Analyzing Response Profiles Using PROC MIXED in SAS](#)
- [5.10 Further Reading](#)

## **Chapter 6: Modeling the Mean: Parametric Curves**

- [6.1 Introduction](#)
- [6.2 Polynomial Trends in Time](#)
- [6.3 Linear Splines](#)
- [6.4 General Linear Model Formulation](#)
- [6.5 Case Studies](#)
- [6.6 Computing: Fitting Parametric Curves Using PROC MIXED in SAS](#)
- [6.7 Further Reading](#)

## **Chapter 7: Modeling the Covariance**

- [7.1 Introduction](#)
- [7.2 Implications of Correlation among Longitudinal Data](#)
- [7.3 Unstructured Covariance](#)
- [7.4 Covariance Pattern Models](#)
- [7.5 Choice among Covariance Pattern Models](#)
- [7.6 Case Study](#)
- [7.7 Discussion: Strengths and Weaknesses of Covariance Pattern Models](#)
- [7.8 Computing: Fitting Covariance Pattern Models Using PROC MIXED in SAS](#)

## [7.9 Further Reading](#)

# [\*\*Chapter 8: Linear Mixed Effects Models\*\*](#)

- [8.1 Introduction](#)
- [8.2 Linear Mixed Effects Models](#)
- [8.3 Random Effects Covariance Structure](#)
- [8.4 Two-Stage Random Effects Formulation](#)
- [8.5 Choice among Random Effects Covariance Models](#)
- [8.6 Prediction of Random Effects](#)
- [8.7 Prediction and Shrinkage](#)
- [8.8 Case Studies](#)
- [8.9 Computing: Fitting Linear Mixed Effects Models Using PROC MIXED in SAS](#)
- [8.10 Further Reading](#)

# [\*\*Chapter 9: Fixed Effects versus Random Effects Models\*\*](#)

- [9.1 Introduction](#)
- [9.2 Linear Fixed Effects Models](#)
- [9.3 Fixed Effects versus Random Effects: Bias-Variance Trade-off](#)
- [9.4 Resolving the Dilemma of Choosing Between Fixed and Random Effects Models](#)
- [9.5 Longitudinal and Cross-sectional Information](#)
- [9.6 Case Study](#)
- [9.7 Computing: Fitting Linear Fixed Effects Models Using PROC GLM in SAS](#)
- [9.8 Computing: Decomposition of Between-Subject and Within-Subject Effects Using PROC MIXED in SAS](#)
- [9.9 Further Reading](#)

# [\*\*Chapter 10: Residual Analyses and Diagnostics\*\*](#)

- [10.1 Introduction](#)
- [10.2 Residuals](#)
- [10.3 Transformed Residuals](#)
- [10.4 Aggregating Residuals](#)
- [10.5 Semi-Variogram](#)
- [10.6 Case Study](#)
- [10.7 Summary](#)
- [10.8 Further Reading](#)

# [\*\*Part III: Generalized Linear Models for Longitudinal Data\*\*](#)

## [\*\*Chapter 11: Review of Generalized Linear Models\*\*](#)

- [11.1 Introduction](#)
- [11.2 Salient Features of Generalized Linear Models](#)
- [11.3 Illustrative Examples](#)
- [11.4 Ordinal Regression Models](#)
- [11.5 Overdispersion](#)
- [11.6 Computing: Fitting Generalized Linear Models Using PROC GENMOD in SAS](#)
- [11.7 Overview of Generalized Linear Models](#)
- [11.8 Further Reading](#)

## **Chapter 12: Marginal Models: Introduction and Overview**

- [12.1 Introduction](#)
- [12.2 Marginal Models for Longitudinal Data](#)
- [12.3 Illustrative Examples of Marginal Models](#)
- [12.4 Distributional Assumptions for Marginal Models](#)
- [12.5 Further Reading](#)

## **Chapter 13: Marginal Models: Generalized Estimating Equations (GEE)**

- [13.1 Introduction](#)
- [13.2 Estimation of Marginal Models: Generalized Estimating Equations](#)
- [13.3 Residual Analyses and Diagnostics](#)
- [13.4 Case Studies](#)
- [13.5 Marginal Models and Time-Varying Covariates](#)
- [13.6 Computing: Generalized Estimating Equations Using PROC GENMOD in SAS](#)
- [13.7 Further Reading](#)

## **Chapter 14: Generalized Linear Mixed Effects Models**

- [14.1 Introduction](#)
- [14.2 Incorporating Random Effects in Generalized Linear Models](#)
- [14.3 Interpretation of Regression Parameters](#)
- [14.4 Overdispersion](#)
- [14.5 Estimation and Inference](#)
- [14.6 A Note on Conditional Maximum Likelihood](#)
- [14.7 Case Studies](#)
- [14.8 Computing: Fitting Generalized Linear Mixed Models Using PROC GLIMMIX in SAS](#)
- [14.9 Further Reading](#)

## **Chapter 15: Generalized Linear Mixed Effects Models: Approximate Methods of Estimation**

- [15.1 Introduction](#)

- [15.2 Penalized Quasi-Likelihood](#)
- [15.3 Marginal Quasi-Likelihood](#)
- [15.4 Cautionary Remarks on the Use of PQL and MQL](#)
- [15.5 Case Studies](#)
- [15.6 Computing: Fitting GLMMs Using PROC GLIMMIX in SAS](#)
- [15.7 Basis of PQL and MQL Approximations](#)
- [15.8 Further Reading](#)

## **Chapter 16: Contrasting Marginal and Mixed Effects Models**

- [16.1 Introduction](#)
- [16.2 Linear Models: A Special Case](#)
- [16.3 Generalized Linear Models](#)
- [16.4 Simple Numerical Illustration](#)
- [16.5 Case Study](#)
- [16.6 Conclusion](#)
- [16.7 Further Reading](#)

## **Part IV: Missing Data and Dropout**

### **Chapter 17: Missing Data and Dropout: Overview of Concepts and Methods**

- [17.1 Introduction](#)
- [17.2 Hierarchy of Missing Data Mechanisms](#)
- [17.3 Implications for Longitudinal Analysis](#)
- [17.4 Dropout](#)
- [17.5 Common Approaches for Handling Dropout](#)
- [17.6 Bias of Last Value Carried Forward Imputation](#)
- [17.7 Further Reading](#)

### **Chapter 18: Missing Data and Dropout: Multiple Imputation and Weighting Methods**

- [18.1 Introduction](#)
- [18.2 Multiple Imputation](#)
- [18.3 Inverse Probability Weighted Methods](#)
- [18.4 Case Studies](#)
- [18.5 “Sandwich” Variance Estimator Adjusting for Estimation of Weights](#)
- [18.6 Computing: Multiple Imputation Using PROC MI in SAS](#)
- [18.7 Computing: Inverse Probability Weighted \(IPW\) Methods in SAS](#)
- [18.8 Further Reading](#)

## **Part V Advanced Topics for Longitudinal and Clustered**

# Data

## Chapter 19: Smoothing Longitudinal Data: Semiparametric Regression Models

- [19.1 Introduction](#)
- [19.2 Penalized Splines for a Univariate Response](#)
- [19.3 Case Study](#)
- [19.4 Penalized Splines for Longitudinal Data](#)
- [19.5 Case Study](#)
- [19.6 Fitting Smooth Curves to Individual Longitudinal Data](#)
- [19.7 Case Study](#)
- [19.8 Computing: Fitting Smooth Curves Using PROC MIXED in SAS](#)
- [19.9 Further Reading](#)

## Chapter 20: Sample Size and Power

- [20.1 Introduction](#)
- [20.2 Sample Size for a Univariate Continuous Response](#)
- [20.3 Sample Size for a Longitudinal Continuous Response](#)
- [20.4 Sample Size for a Longitudinal Binary Response](#)
- [20.5 Summary](#)
- [20.6 Computing: Sample Size Calculation Using Pseudo-Data](#)
- [20.7 Further Reading](#)

## Chapter 21: Repeated Measures and Related Designs

- [21.1 Introduction](#)
- [21.2 Repeated Measures Designs](#)
- [21.3 Multiple Source Data](#)
- [21.4 Case Study 1: Repeated Measures Experiment](#)
- [21.5 Case Study 2: Multiple Source Data](#)
- [21.6 Summary](#)
- [21.7 Further Reading](#)

## Chapter 22: Multilevel Models

- [22.1 Introduction](#)
- [22.2 Multilevel Data](#)
- [22.3 Multilevel Linear Models](#)
- [22.4 Multilevel Generalized Linear Models](#)
- [22.5 Summary](#)
- [22.6 Further Reading](#)

## Appendix A Gentle Introduction to Vectors and Matrices

[Appendix B Properties of Expectations and Variances](#)

[Appendix C Critical Points for a 50:50 Mixture of Chi-Squared Distributions](#)

[References](#)

[Index](#)

# **Applied Longitudinal Analysis**

# WILEY SERIES IN PROBABILITY AND STATISTICS

Established by WALTER A. SHEWHART and SAMUEL S. WILKS

Editors: *David J. Balding, Noel A. C. Cressie, Garrett M. Fitzmaurice, Harvey Goldstein, Iain M. Johnstone, Geert Molenberghs, David W. Scott, Adrian F. M. Smith, Ruey S. Tsay, Sanford Weisberg* Editors Emeriti: *Vic Barnett, J. Stuart Hunter, Joseph B. Kadane, Jozef L. Teugels*

The *Wiley Series in Probability and Statistics* is well established and authoritative. It covers many topics of current research interest in both pure and applied statistics and probability theory. Written by leading statisticians and institutions, the titles span both state-of-the-art developments in the field and classical methods.

Reflecting the wide range of current research in statistics, the series encompasses applied, methodological and theoretical statistics, ranging from applications and new techniques made possible by advances in computerized practice to rigorous treatment of theoretical approaches.

This series provides essential and invaluable reading for all statisticians, whether in academia, industry, government, or research.

<sup>†</sup> ABRAHAM and LEDOLTER • Statistical Methods for Forecasting

AGRESTI • Analysis of Ordinal Categorical Data, *Second Edition*

AGRESTI • An Introduction to Categorical Data Analysis, *Second Edition*

AGRESTI • Categorical Data Analysis, *Second Edition*

ALTMAN, GILL, and McDONALD • Numerical Issues in Statistical Computing for the Social Scientist

AMARATUNGA and CABRERA • Exploration and Analysis of DNA Microarray and Protein Array Data

ANDĚL • Mathematics of Chance

ANDERSON • An Introduction to Multivariate Statistical Analysis, *Third Edition*

\* ANDERSON • The Statistical Analysis of Time Series

ANDERSON, AUQUIER, HAUCK, OAKES, VANDAELE, and WEISBERG • Statistical Methods for Comparative Studies

ANDERSON and LOYNES • The Teaching of Practical Statistics

ARMITAGE and DAVID (editors) • Advances in Biometry

ARNOLD, BALAKRISHNAN, and NAGARAJA • Records

\* ARTHANARI and DODGE • Mathematical Programming in Statistics

\* BAILEY • The Elements of Stochastic Processes with Applications to the Natural Sciences

BALAKRISHNAN and KOUTRAS • Runs and Scans with Applications

BALAKRISHNAN and NG • Precedence-Type Tests and Applications

BARNETT • Comparative Statistical Inference, *Third Edition*

BARNETT • Environmental Statistics

BARNETT and LEWIS • Outliers in Statistical Data, *Third Edition*

BARTOSZYNSKI and NIEWIADOMSKA-BUGAJ • Probability and Statistical Inference

BASILEVSKY • Statistical Factor Analysis and Related Methods: Theory and Applications

BASU and RIGDON • Statistical Methods for the Reliability of Repairable Systems

BATES and WATTS • Nonlinear Regression Analysis and Its Applications

BECHHOFER, SANTNER, and GOLDSMAN • Design and Analysis of Experiments for Statistical Selection, Screening, and Multiple Comparisons

BELSLEY • Conditioning Diagnostics: Collinearity and Weak Data in Regression

<sup>†</sup> BELSLEY, KUH, and WELSCH • Regression Diagnostics: Identifying Influential Data and Sources of Collinearity

BENDAT and PIERSOL • Random Data: Analysis and Measurement Procedures, *Fourth Edition*

BERRY, CHALONER, and GEWEKE • Bayesian Analysis in Statistics and Econometrics: Essays in

Honor of Arnold Zellner  
BERNARDO and SMITH • Bayesian Theory  
BHAT and MILLER • Elements of Applied Stochastic Processes, *Third Edition*  
BHATTACHARYA and WAYMIRE • Stochastic Processes with Applications  
BILLINGSLEY • Convergence of Probability Measures, *Second Edition*  
BILLINGSLEY • Probability and Measure, *Third Edition*  
BIRKES and DODGE • Alternative Methods of Regression  
BISGAARD and KULAHCI • Time Series Analysis and Forecasting by Example  
BISWAS, DATTA, FINE, and SEGAL • Statistical Advances in the Biomedical Sciences: Clinical Trials, Epidemiology, Survival Analysis, and Bioinformatics  
BLISCHKE AND MURTHY (editors) • Case Studies in Reliability and Maintenance  
BLISCHKE AND MURTHY • Reliability: Modeling, Prediction, and Optimization  
BLOOMFIELD • Fourier Analysis of Time Series: An Introduction, *Second Edition*  
BOLLEN • Structural Equations with Latent Variables  
BOLLEN and CURRAN • Latent Curve Models: A Structural Equation Perspective  
BOROVKOV • Ergodicity and Stability of Stochastic Processes  
BOULEAU • Numerical Methods for Stochastic Processes  
BOX • Bayesian Inference in Statistical Analysis  
BOX • R. A. Fisher, the Life of a Scientist  
BOX and DRAPER • Response Surfaces, Mixtures, and Ridge Analyses, *Second Edition*  
\* BOX and DRAPER • Evolutionary Operation: A Statistical Method for Process Improvement  
BOX and FRIENDS • Improving Almost Anything, *Revised Edition*  
BOX, HUNTER, and HUNTER • Statistics for Experimenters: Design, Innovation, and Discovery, *Second Edition*  
BOX, JENKINS, and REINSEL • Time Series Analysis: Forecasting and Control, *Fourth Edition*  
BOX, LUCEÑO, and PANIAGUA-QUIÑONES • Statistical Control by Monitoring and Adjustment, *Second Edition*  
BRANDIMARTE • Numerical Methods in Finance: A MATLAB-Based Introduction  
† BROWN and HOLLANDER • Statistics: A Biomedical Introduction  
BRUNNER, DOMHOF, and LANGER • Nonparametric Analysis of Longitudinal Data in Factorial Experiments  
BUCKLEW • Large Deviation Techniques in Decision, Simulation, and Estimation  
CAIROLI and DALANG • Sequential Stochastic Optimization  
CASTILLO, HADI, BALAKRISHNAN, and SARABIA • Extreme Value and Related Models with Applications in Engineering and Science  
CHAN • Time Series: Applications to Finance with R and S-Plus®, *Second Edition*  
CHARALAMBIDES • Combinatorial Methods in Discrete Distributions  
CHATTERJEE and HADI • Regression Analysis by Example, *Fourth Edition*  
CHATTERJEE and HADI • Sensitivity Analysis in Linear Regression  
CHERNICK • Bootstrap Methods: A Guide for Practitioners and Researchers, *Second Edition*  
CHERNICK and FRIIS • Introductory Biostatistics for the Health Sciences  
CHILÈS and DELFINER • Geostatistics: Modeling Spatial Uncertainty  
CHOW and LIU • Design and Analysis of Clinical Trials: Concepts and Methodologies, *Second Edition*  
CLARKE • Linear Models: The Theory and Application of Analysis of Variance  
CLARKE and DISNEY • Probability and Random Processes: A First Course with Applications, *Second Edition*  
\* COCHRAN and COX • Experimental Designs, *Second Edition*  
COLLINS and LANZA • Latent Class and Latent Transition Analysis: With Applications in the

Social, Behavioral, and Health Sciences

CONGDON • Applied Bayesian Modelling

CONGDON • Bayesian Models for Categorical Data

CONGDON • Bayesian Statistical Modelling

CONOVER • Practical Nonparametric Statistics, *Third Edition*

COOK • Regression Graphics

COOK and WEISBERG • Applied Regression Including Computing and Graphics

COOK and WEISBERG • An Introduction to Regression Graphics

CORNELL • Experiments with Mixtures, Designs, Models, and the Analysis of Mixture Data, *Third Edition*

COVER and THOMAS • Elements of Information Theory

COX • A Handbook of Introductory Statistical Methods

\* COX • Planning of Experiments

CRESSIE • Statistics for Spatial Data, *Revised Edition*

CRESSIE and WIKLE • Statistics for Spatio-Temporal Data

CSÖRGÖ and HORVÁTH • Limit Theorems in Change Point Analysis

DANIEL • Applications of Statistics to Industrial Experimentation

DANIEL • Biostatistics: A Foundation for Analysis in the Health Sciences, *Eighth Edition*

\* DANIEL • Fitting Equations to Data: Computer Analysis of Multifactor Data, *Second Edition*

DASU and JOHNSON • Exploratory Data Mining and Data Cleaning

DAVID and NAGARAJA • Order Statistics, *Third Edition*

\* DEGROOT, FIENBERG, and KADANE • Statistics and the Law

DEL CASTILLO • Statistical Process Adjustment for Quality Control

DEMARIS • Regression with Social Data: Modeling Continuous and Limited Response Variables

DEMIDENKO • Mixed Models: Theory and Applications

DENISON, HOLMES, MALLICK and SMITH • Bayesian Methods for Nonlinear Classification and Regression

DETTE and STUDDEN • The Theory of Canonical Moments with Applications in Statistics, Probability, and Analysis

DEY and MUKERJEE • Fractional Factorial Plans

DILLON and GOLDSTEIN • Multivariate Analysis: Methods and Applications

DODGE • Alternative Methods of Regression

\* DODGE and ROMIG • Sampling Inspection Tables, *Second Edition*

\* DOOB • Stochastic Processes

DOWDY, WEARDEN, and CHILKO • Statistics for Research, *Third Edition*

DRAPER and SMITH • Applied Regression Analysis, *Third Edition*

DRYDEN and MARDIA • Statistical Shape Analysis

DUDEWICZ and MISHRA • Modern Mathematical Statistics

DUNN and CLARK • Basic Statistics: A Primer for the Biomedical Sciences, *Third Edition*

DUPUIS and ELLIS • A Weak Convergence Approach to the Theory of Large Deviations

EDLER and KITSOS • Recent Advances in Quantitative Methods in Cancer and Human Health Risk Assessment

\* ELANDT-JOHNSON and JOHNSON • Survival Models and Data Analysis

ENDERS • Applied Econometric Time Series

† ETHIER and KURTZ • Markov Processes: Characterization and Convergence

EVANS, HASTINGS, and PEACOCK • Statistical Distributions, *Third Edition*

FELLER • An Introduction to Probability Theory and Its Applications, Volume I, *Third Edition*, Revised; Volume II, *Second Edition*

FISHER and VAN BELLE • Biostatistics: A Methodology for the Health Sciences

- FITZMAURICE, LAIRD, and WARE • Applied Longitudinal Analysis, *Second Edition*
- \* FLEISS • The Design and Analysis of Clinical Experiments
- FLEISS • Statistical Methods for Rates and Proportions, *Third Edition*
- † FLEMING and HARRINGTON • Counting Processes and Survival Analysis
- FUJIKOSHI, ULYANOV, and SHIMIZU • Multivariate Statistics: High-Dimensional and Large-Sample Approximations
- FULLER • Introduction to Statistical Time Series, *Second Edition*
- † FULLER • Measurement Error Models
- GALLANT • Nonlinear Statistical Models
- GEISSER • Modes of Parametric Statistical Inference
- GELMAN and MENG • Applied Bayesian Modeling and Causal Inference from Incomplete-Data Perspectives
- GEWEKE • Contemporary Bayesian Econometrics and Statistics
- GHOSH, MUKHOPADHYAY, and SEN • Sequential Estimation
- GIESBRECHT and GUMPERTZ • Planning, Construction, and Statistical Analysis of Comparative Experiments
- GIFI • Nonlinear Multivariate Analysis
- GIVENS and HOETING • Computational Statistics
- GLASSERMAN and YAO • Monotone Structure in Discrete-Event Systems
- GNANADESIKAN • Methods for Statistical Data Analysis of Multivariate Observations, *Second Edition*
- GOLDSTEIN and LEWIS • Assessment: Problems, Development, and Statistical Issues
- GREENWOOD and NIKULIN • A Guide to Chi-Squared Testing
- GROSS, SHORTLE, THOMPSON, and HARRIS • Fundamentals of Queueing Theory, *Fourth Edition*
- GROSS, SHORTLE, THOMPSON, and HARRIS • Solutions Manual to Accompany Fundamentals of Queueing Theory, *Fourth Edition*
- \* HAHN and SHAPIRO • Statistical Models in Engineering
- HAHN and MEEKER • Statistical Intervals: A Guide for Practitioners
- HALD • A History of Probability and Statistics and their Applications Before 1750
- HALD • A History of Mathematical Statistics from 1750 to 1930
- † HAMPEL • Robust Statistics: The Approach Based on Influence Functions
- HANNAN and DEISTLER • The Statistical Theory of Linear Systems
- HARMAN and KULKARNI • An Elementary Introduction to Statistical Learning Theory
- HARTUNG, KNAPP, and SINHA • Statistical Meta-Analysis with Applications
- HEIBERGER • Computation for the Analysis of Designed Experiments
- HEDAYAT and SINHA • Design and Inference in Finite Population Sampling
- HEDEKER and GIBBONS • Longitudinal Data Analysis
- HELLER • MACSYMA for Statisticians
- HINKELMANN and KEMPTHORNE • Design and Analysis of Experiments, Volume 1: Introduction to Experimental Design, *Second Edition*
- HINKELMANN and KEMPTHORNE • Design and Analysis of Experiments, Volume 2: Advanced Experimental Design
- HOAGLIN, MOSTELLER, and TUKEY • Fundamentals of Exploratory Analysis of Variance
- \* HOAGLIN, MOSTELLER, and TUKEY • Exploring Data Tables, Trends and Shapes
- \* HOAGLIN, MOSTELLER, and TUKEY • Understanding Robust and Exploratory Data Analysis
- HOCHBERG and TAMHANE • Multiple Comparison Procedures
- HOCKING • Methods and Applications of Linear Models: Regression and the Analysis of Variance, *Second Edition*

- HOEL • Introduction to Mathematical Statistics, *Fifth Edition*  
HOGG and KLUGMAN • Loss Distributions  
HOLLANDER and WOLFE • Nonparametric Statistical Methods, *Second Edition*  
HOSMER and LEMESHOW • Applied Logistic Regression, *Second Edition*  
HOSMER, LEMESHOW, and MAY • Applied Survival Analysis: Regression Modeling of Time-to-Event Data, *Second Edition*  
HUBER • Data Analysis: What Can Be Learned From the Past 50 Years  
<sup>†</sup> HUBER and RONCHETTI • Robust Statistics, *Second Edition*  
HUBERTY • Applied Discriminant Analysis  
HUBERTY • Applied Discriminant Analysis *Second Edition*  
HUNT and KENNEDY • Financial Derivatives in Theory and Practice, *Revised Edition*  
HURD and MIAMEE • Periodically Correlated Random Sequences: Spectral Theory and Practice  
HUSKOVA, BERAN, and DUPAC • Collected Works of Jaroslav Hajek—with Commentary  
HUZURBAZAR • Flowgraph Models for Multistate Time-to-Event Data  
IMAN and CONOVER • A Modern Approach to Statistics  
<sup>\*</sup> JACKSON • A User's Guide to Principle Components  
JOHN • Statistical Methods in Engineering and Quality Assurance  
JOHNSON • Multivariate Statistical Simulation  
JOHNSON and BALAKRISHNAN • Advances in the Theory and Practice of Statistics: A Volume in Honor of Samuel Kotz  
JOHNSON and BHATTACHARYYA • Statistics: Principles and Methods, *Fifth Edition*  
JOHNSON and KOTZ • Distributions in Statistics  
JOHNSON and KOTZ (editors) • Leading Personalities in Statistical Sciences: From the Seventeenth Century to the Present  
JOHNSON, KOTZ, and BALAKRISHNAN • Continuous Univariate Distributions, Volume 1, *Second Edition*  
JOHNSON, KOTZ, and BALAKRISHNAN • Continuous Univariate Distributions, Volume 2, *Second Edition*  
JOHNSON, KOTZ, and BALAKRISHNAN • Discrete Multivariate Distributions  
JOHNSON, KEMP, and KOTZ • Univariate Discrete Distributions, *Third Edition*  
JUDGE, GRIFFITHS, HILL, LÜTKEPOHL, and LEE • The Theory and Practice of Econometrics, *Second Edition*  
JUREČKOVÁ and SEN • Robust Statistical Procedures: Asymptotics and Interrelations  
JUREK and MASON • Operator-Limit Distributions in Probability Theory  
KADANE • Bayesian Methods and Ethics in a Clinical Trial Design  
KADANE AND SCHUM • A Probabilistic Analysis of the Sacco and Vanzetti Evidence  
KALBFLEISCH and PRENTICE • The Statistical Analysis of Failure Time Data, *Second Edition*  
KARIYA and KURATA • Generalized Least Squares  
KASS and VOS • Geometrical Foundations of Asymptotic Inference  
<sup>†</sup> KAUFMAN and ROUSSEEUW • Finding Groups in Data: An Introduction to Cluster Analysis  
KEDEM and FOKIANOS • Regression Models for Time Series Analysis  
KENDALL, BARDEEN, CARNE, and LE • Shape and Shape Theory  
KHURI • Advanced Calculus with Applications in Statistics, *Second Edition*  
KHURI, MATHEW, and SINHA • Statistical Tests for Mixed Linear Models  
KLEIBER and KOTZ • Statistical Size Distributions in Economics and Actuarial Sciences  
KLEMELÄ • Smoothing of Multivariate Data: Density Estimation and Visualization  
KLUGMAN, PANJER, and WILLMOT • Loss Models: From Data to Decisions, *Third Edition*  
KLUGMAN, PANJER, and WILLMOT • Solutions Manual to Accompany Loss Models: From Data to Decisions, *Third Edition*

- KOTZ, BALAKRISHNAN, and JOHNSON • Continuous Multivariate Distributions, Volume 1, *Second Edition*
- KOVALENKO, KUZNETZOV, and PEGG • Mathematical Theory of Reliability of Time-Dependent Systems with Practical Applications
- KOWALSKI and TU • Modern Applied U-Statistics
- KRISHNAMOORTHY and MATHEW • Statistical Tolerance Regions: Theory, Applications, and Computation
- KROESE, TAIMRE, and BOTEV • Handbook of Monte Carlo Methods
- KROONENBERG • Applied Multiway Data Analysis
- KVAM and VIDAKOVIC • Nonparametric Statistics with Applications to Science and Engineering
- LACHIN • Biostatistical Methods: The Assessment of Relative Risks, *Second Edition*
- LAD • Operational Subjective Statistical Methods: A Mathematical, Philosophical, and Historical Introduction
- LAMPERTI • Probability: A Survey of the Mathematical Theory, *Second Edition*
- LANGE, RYAN, BILLARD, BRILLINGER, CONQUEST, and GREENHOUSE • Case Studies in Biometry
- LARSON • Introduction to Probability Theory and Statistical Inference, *Third Edition*
- LAWLESS • Statistical Models and Methods for Lifetime Data, *Second Edition*
- LAWSON • Statistical Methods in Spatial Epidemiology
- LE • Applied Categorical Data Analysis
- LE • Applied Survival Analysis
- LEE and WANG • Statistical Methods for Survival Data Analysis, *Third Edition*
- LEPAGE and BILLARD • Exploring the Limits of Bootstrap
- LEYLAND and GOLDSTEIN (editors) • Multilevel Modelling of Health Statistics
- LIAO • Statistical Group Comparison
- LINDVALL • Lectures on the Coupling Method
- LIN • Introductory Stochastic Analysis for Finance and Insurance
- LINHART and ZUCCHINI • Model Selection
- LITTLE and RUBIN • Statistical Analysis with Missing Data, *Second Edition*
- LLOYD • The Statistical Analysis of Categorical Data
- LOWEN and TEICH • Fractal-Based Point Processes
- MAGNUS and NEUDECKER • Matrix Differential Calculus with Applications in Statistics and Econometrics, *Revised Edition*
- MALLER and ZHOU • Survival Analysis with Long Term Survivors
- MALLOWS • Design, Data, and Analysis by Some Friends of Cuthbert Daniel
- MANN, SCHAFER, and SINGPURWALLA • Methods for Statistical Analysis of Reliability and Life Data
- MANTON, WOODBURY, and TOLLEY • Statistical Applications Using Fuzzy Sets
- MARCHETTE • Random Graphs for Statistical Pattern Recognition
- MARDIA and JUPP • Directional Statistics
- MASON, GUNST, and HESS • Statistical Design and Analysis of Experiments with Applications to Engineering and Science, *Second Edition*
- McCULLOCH, SEARLE, and NEUHAUS • Generalized, Linear, and Mixed Models, *Second Edition*
- McFADDEN • Management of Data in Clinical Trials, *Second Edition*
- \* McLACHLAN • Discriminant Analysis and Statistical Pattern Recognition
- McLACHLAN, DO, and AMBROISE • Analyzing Microarray Gene Expression Data
- McLACHLAN and KRISHNAN • The EM Algorithm and Extensions, *Second Edition*
- McLACHLAN and PEEL • Finite Mixture Models
- McNEIL • Epidemiological Research Methods

- MEEKER and ESCOBAR • Statistical Methods for Reliability Data
- MEERSCHAERT and SCHEFFLER • Limit Distributions for Sums of Independent Random Vectors: Heavy Tails in Theory and Practice
- MICKEY, DUNN, and CLARK • Applied Statistics: Analysis of Variance and Regression, *Third Edition*
- \* MILLER • Survival Analysis, *Second Edition*
- MONTGOMERY, JENNINGS, and KULAHCI • Introduction to Time Series Analysis and Forecasting
- MONTGOMERY, PECK, and VINING • Introduction to Linear Regression Analysis, *Fourth Edition*
- MORGENTHALER and TUKEY • Configural Polysampling: A Route to Practical Robustness
- MUIRHEAD • Aspects of Multivariate Statistical Theory
- MULLER and STOYAN • Comparison Methods for Stochastic Models and Risks
- MURRAY • X-STAT 2.0 Statistical Experimentation, Design Data Analysis, and Nonlinear Optimization
- MURTHY, XIE, and JIANG • Weibull Models
- MYERS, MONTGOMERY, and ANDERSON-COOK • Response Surface Methodology: Process and Product Optimization Using Designed Experiments, *Third Edition*
- MYERS, MONTGOMERY, VINING, and ROBINSON • Generalized Linear Models. With Applications in Engineering and the Sciences, *Second Edition*
- \* NELSON • Accelerated Testing, Statistical Models, Test Plans, and Data Analyses
- † NELSON • Applied Life Data Analysis
- NEWMAN • Biostatistical Methods in Epidemiology
- OCHI • Applied Probability and Stochastic Processes in Engineering and Physical Sciences
- OKABE, BOOTS, SUGIHARA, and CHIU • Spatial Tesselations: Concepts and Applications of Voronoi Diagrams, *Second Edition*
- OLIVER and SMITH • Influence Diagrams, Belief Nets and Decision Analysis
- PALTA • Quantitative Methods in Population Health: Extensions of Ordinary Regressions
- PANJER • Operational Risk: Modeling and Analytics
- PANKRATZ • Forecasting with Dynamic Regression Models
- PANKRATZ • Forecasting with Univariate Box-Jenkins Models: Concepts and Cases
- \* PARZEN • Modern Probability Theory and Its Applications
- PEÑA, TIAO, and TSAY • A Course in Time Series Analysis
- PIANTADOSI • Clinical Trials: A Methodologic Perspective
- PORT • Theoretical Probability for Applications
- POURAHMADI • Foundations of Time Series Analysis and Prediction Theory POWELL • Approximate Dynamic Programming: Solving the Curses of Dimensionality
- PRESS • Bayesian Statistics: Principles, Models, and Applications
- PRESS • Subjective and Objective Bayesian Statistics, *Second Edition*
- PRESS and TANUR • The Subjectivity of Scientists and the Bayesian Approach
- PUKELSHEIM • Optimal Experimental Design
- PURI, VILAPLANA, and WERTZ • New Perspectives in Theoretical and Applied Statistics
- † PUTERMAN • Markov Decision Processes: Discrete Stochastic Dynamic Programming
- QIU • Image Processing and Jump Regression Analysis
- \* RAO • Linear Statistical Inference and Its Applications, *Second Edition*
- RAUSAND and HØYLAND • System Reliability Theory: Models, Statistical Methods, and Applications, *Second Edition*
- RENCHER • Linear Models in Statistics
- RENCHER • Methods of Multivariate Analysis, *Second Edition*
- RENCHER • Multivariate Statistical Inference with Applications

- \* RIPLEY • Spatial Statistics
- \* RIPLEY • Stochastic Simulation
- ROBINSON • Practical Strategies for Experimenting
- ROHATGI and SALEH • An Introduction to Probability and Statistics, *Second Edition*
- ROLSKI, SCHMIDLI, SCHMIDT, and TEUGELS • Stochastic Processes for Insurance and Finance
- ROSENBERGER and LACHIN • Randomization in Clinical Trials: Theory and Practice
- ROSS • Introduction to Probability and Statistics for Engineers and Scientists
- ROSSI, ALLENBY, and McCULLOCH • Bayesian Statistics and Marketing
- † ROUSSEEUW and LEROY • Robust Regression and Outlier Detection
- \* RUBIN • Multiple Imputation for Nonresponse in Surveys
- RUBINSTEIN and KROESE • Simulation and the Monte Carlo Method, *Second Edition*
- RUBINSTEIN and MELAMED • Modern Simulation and Modeling
- RYAN • Modern Engineering Statistics
- RYAN • Modern Experimental Design
- RYAN • Modern Regression Methods, *Second Edition*
- RYAN • Statistical Methods for Quality Improvement, *Second Edition*
- SALEH • Theory of Preliminary Test and Stein-Type Estimation with Applications
- \* SCHEFFE • The Analysis of Variance
- SCHIMEK • Smoothing and Regression: Approaches, Computation, and Application
- SCHOTT • Matrix Analysis for Statistics, *Second Edition*
- SCHOUTENS • Levy Processes in Finance: Pricing Financial Derivatives
- SCHUSS • Theory and Applications of Stochastic Differential Equations
- SCOTT • Multivariate Density Estimation: Theory, Practice, and Visualization
- † SEARLE • Linear Models for Unbalanced Data
- † SEARLE • Matrix Algebra Useful for Statistics
- † SEARLE, CASELLA, and McCULLOCH • Variance Components
- SEARLE and WILLETT • Matrix Algebra for Applied Economics
- SEBER • A Matrix Handbook For Statisticians
- † SEBER • Multivariate Observations
- SEBER and LEE • Linear Regression Analysis, *Second Edition*
- † SEBER and WILD • Nonlinear Regression
- SENNOTT • Stochastic Dynamic Programming and the Control of Queueing Systems
- \* SERFLING • Approximation Theorems of Mathematical Statistics
- SHAFER and VOVK • Probability and Finance: It's Only a Game!
- SILVAPULLE and SEN • Constrained Statistical Inference: Inequality, Order, and Shape Restrictions
- SMALL and McLEISH • Hilbert Space Methods in Probability and Statistical Inference
- SRIVASTAVA • Methods of Multivariate Statistics
- STAPLETON • Linear Statistical Models, *Second Edition*
- STAPLETON • Models for Probability and Statistical Inference: Theory and Applications
- STAUDTE and SHEATHER • Robust Estimation and Testing
- STOYAN, KENDALL, and MECKE • Stochastic Geometry and Its Applications, *Second Edition*
- STOYAN and STOYAN • Fractals, Random Shapes and Point Fields: Methods of Geometrical Statistics
- STREET and BURGESS • The Construction of Optimal Stated Choice Experiments: Theory and Methods
- STYAN • The Collected Papers of T. W. Anderson: 1943–1985
- SUTTON, ABRAMS, JONES, SHELDON, and SONG • Methods for Meta-Analysis in Medical Research

- TAKEZAWA • Introduction to Nonparametric Regression
- TAMHANE • Statistical Analysis of Designed Experiments: Theory and Applications
- TANAKA • Time Series Analysis: Nonstationary and Noninvertible Distribution Theory
- THOMPSON • Empirical Model Building
- THOMPSON • Sampling, *Second Edition*
- THOMPSON • Simulation: A Modeler's Approach
- THOMPSON and SEBER • Adaptive Sampling
- THOMPSON, WILLIAMS, and FINDLAY • Models for Investors in Real World Markets
- TCIAO, BISGAARD, HILL, PEÑA, and STIGLER (editors) • Box on Quality and Discovery: with Design, Control, and Robustness
- TIERNEY • LISP-STAT: An Object-Oriented Environment for Statistical Computing and Dynamic Graphics
- TSAY • Analysis of Financial Time Series, *Third Edition*
- UPTON and FINGLETON • Spatial Data Analysis by Example, Volume II: Categorical and Directional Data
- † VAN BELLE • Statistical Rules of Thumb, *Second Edition*
- VAN BELLE, FISHER, HEAGERTY, and LUMLEY • Biostatistics: A Methodology for the Health Sciences, *Second Edition*
- VESTRUP • The Theory of Measures and Integration
- VIDAKOVIC • Statistical Modeling by Wavelets
- VINOD and REAGLE • Preparing for the Worst: Incorporating Downside Risk in Stock Market Investments
- WALLER and GOTWAY • Applied Spatial Statistics for Public Health Data
- WEERAHANDI • Generalized Inference in Repeated Measures: Exact Methods in MANOVA and Mixed Models
- WEISBERG • Applied Linear Regression, *Third Edition*
- WEISBERG • Bias and Causation: Models and Judgment for Valid Comparisons
- WELSH • Aspects of Statistical Inference
- WESTFALL and YOUNG • Resampling-Based Multiple Testing: Examples and Methods for *p*-Value Adjustment
- WHITTAKER • Graphical Models in Applied Multivariate Statistics
- WINKER • Optimization Heuristics in Economics: Applications of Threshold Accepting
- WONNACOTT and WONNACOTT • Econometrics, *Second Edition*
- WOODING • Planning Pharmaceutical Clinical Trials: Basic Statistical Principles
- WOODWORTH • Biostatistics: A Bayesian Introduction
- WOOLSON and CLARKE • Statistical Methods for the Analysis of Biomedical Data, *Second Edition*
- WU and HAMADA • Experiments: Planning, Analysis, and Parameter Design Optimization, *Second Edition*
- WU and ZHANG • Nonparametric Regression Methods for Longitudinal Data Analysis
- YANG • The Construction Theory of Denumerable Markov Processes
- YOUNG, VALERO-MORA, and FRIENDLY • Visual Statistics: Seeing Data with Dynamic Interactive Graphics
- ZACKS • Stage-Wise Adaptive Designs
- ZELTERMAN • Discrete Distributions—Applications in the Health Sciences
- \* ZELLNER • An Introduction to Bayesian Inference in Econometrics
- ZHOU, OBUCHOWSKI, and McCLISH • Statistical Methods in Diagnostic Medicine, *Second Edition*

\* Now available in a lower priced paperback edition in the Wiley Classics Library.

† Now available in a lower priced paperback edition in the Wiley-Interscience Paperback Series.

# **Applied Longitudinal Analysis**

Second Edition

**Garrett M. Fitzmaurice**

**Nan M. Laird**

**James H. Ware**

*Department of Biostatistics*

*Harvard University*

*Boston, MA*



A JOHN WILEY & SONS, INC., PUBLICATION

Copyright © 2011 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, or online at <http://www.wiley.com/go/permission>.

**Limit of Liability/Disclaimer of Warranty:** While the publisher and author have used their best efforts in preparing this book, they make no representation 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 United States at (800) 762-2974, outside the United States 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 formats. For more information about Wiley products, visit our web site at [www.wiley.com](http://www.wiley.com).

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

Fitzmaurice, Garrett M., 1962-

Applied longitudinal analysis / Garrett M. Fitzmaurice, Nan M. Laird, James H. Ware. — 2nd ed.  
p. cm.

ISBN 978-0-470-38027-7 (hardback)

1. Longitudinal method. 2. Regression analysis. 3. Multivariate analysis. 4. Medical statistics. I. Laird, Nan M., 1943- II. Ware, James H., 1941- III. Title.

QA278.F575 2011

519.5'3—dc22

2011012197

*To Laura, Kieran, and Aidan*  
— G.M.F.

*To Joel, Richard, and Lily*  
— N.M.L.

*To Janice, Cameron, and Jake*  
— J.H.W.

# Preface

The first edition of *Applied Longitudinal Analysis* was designed to serve as a textbook for a course on modern statistical methods for longitudinal data analysis, and subsequently, as a reference resource for students and researchers. The book was targeted at a broad audience: graduate students in statistics, statisticians working in the health sciences, pharmaceutical industry, and governmental health-related agencies, as well as researchers and graduate students from a variety of substantive fields. In the seven years that have elapsed since publication of the first edition, *Applied Longitudinal Analysis* has been used extensively in university classrooms throughout the United States and abroad. We are grateful to many colleagues, course instructors, students, and readers who have offered constructive suggestions on how the book could be improved. This feedback has been invaluable and helped shape the content of the second edition.

The feedback we received has encouraged us to retain the general structure and format of the first edition while taking the opportunity to introduce a number of new and important topics. Although there is much new material in this second edition, the principles that guided us in writing the first edition have not changed. Our primary goal is to present a rigorous and comprehensive description of modern statistical methods for the analysis of longitudinal data that is accessible to a wide range of readers. A strong emphasis is placed on the application of these methods to longitudinal data and the interpretation of results. Although the methods are presented in the setting of numerous applications to actual data sets drawn from studies in health-related fields, reflecting our own research interests in the health sciences, they apply equally to other areas of application, for example, education, psychology, and other branches of the behavioral and social sciences.

How does this edition differ from its predecessor? The major changes in this edition have resulted from the addition of six new chapters:

1. A chapter (Chapter 9) on “fixed effects models,” in which subject-specific effects are treated as fixed rather than random, has been added. This chapter complements the existing chapter on mixed effects models (Chapter 8) and includes a discussion of the relative advantages of these two classes of models.
2. In the first edition, a single chapter was devoted to marginal models and generalized estimating equations (GEE) that focused exclusively on binary and count data. We now devote two chapters (Chapters 12 and 13) to marginal models and GEE, with new material on models for ordinal responses, residual diagnostics, and issues that arise when modeling time-varying covariates.
3. A chapter (Chapter 15) on approximate methods for generalized linear mixed effects models discusses penalized quasi-likelihood (PQL) and marginal quasi-likelihood (MQL) methods. We highlight settings where these approximations are unlikely to be accurate and can yield biased estimates of effects.
4. A second chapter (Chapter 18) on missing data and dropout, focusing on multiple imputation and inverse probability weighting (IPW) methods, has been added. To give greater prominence to methods for accounting for missing data and dropout in longitudinal analyses, the two companion chapters (Chapters 17 and 18) now appear before the *Advanced Topics* part of the book.
5. A chapter (Chapter 19) on smoothing longitudinal data has been added to the *Advanced Topics*. This chapter focuses on the connection between penalized splines and linear mixed effects models.
6. A chapter on sample size and power (Chapter 20) has been added to the *Advanced Topics*. This chapter considers issues of sample size, power, number of repeated measurements, and study duration for longitudinal study designs.

In addition the chapter on residual analyses and diagnostics (Chapter 10) has been revised to include material on recently developed model-checking techniques based on cumulative sums of residuals

and the chapters that review generalized linear models (Chapter 11) and generalized linear mixed effects models (Chapter 14) have been updated to include new material on models for ordinal data and on methods for handling overdispersion. Finally, extra problem sets have been added to many of the chapters.

As in the first edition, the prerequisites for a course based on this book are an introductory course in statistics and a strong background in regression analysis. Some previous exposure to generalized linear models (e.g., logistic regression) would be helpful, although these models are reviewed in detail in the text. An understanding of matrix algebra or calculus is not assumed. Although we do not assume a high level of mathematical preparation, we have written this book for the motivated reader who is willing to consider mathematical ideas. The more technical or mathematical sections of the book are signposted with asterisks and may be omitted at first reading without loss of continuity.

The methods described in this book require the use of appropriate statistical software. As before, we include illustrative *SAS* commands for performing the analyses presented throughout the text at the end of many chapters, with basic descriptions of their usage. Because many of the analyses we discuss can be performed using alternative software packages (e.g., *R*, *S-Plus*, *Stata*, and *SPSS*), this book can be supplemented with any one of them. Readers are encouraged to perform and verify the results of analyses using statistical software of their choice. Programming statements and computer output for selected examples, prepared using *SAS*, *Stata*, and *R*, can be downloaded from the website: [www.biostat.harvard.edu/~fitzmaur/ala2e](http://www.biostat.harvard.edu/~fitzmaur/ala2e). Because statistical software is constantly evolving, we will endeavor to update the website as new procedures become available in the major statistical software packages. The thirty-two real data sets used throughout the text and problem sets to illustrate the applications of longitudinal methods also can be downloaded from the website.

We hope this second edition of *Applied Longitudinal Analysis* provides a broader foundation in modern methods for the analysis of longitudinal data and will prove a worthy successor to the first edition. The original impetus for writing this book arose from teaching a graduate-level course on “Applied Longitudinal Analysis” at the Harvard School of Public Health. We are especially grateful to the students who have participated in the course since its inception almost twenty years ago; we have learned much from these extraordinary students. The collection of individuals who gave us useful feedback on the first edition is far too long to list. However, we would like to thank the many friends and colleagues who have helped us with this project. A special word of thanks to Amy Herring and Russell Localio. We thank Amy for her many helpful and constructive suggestions on how the book could be improved. We thank Russell for reading a draft of the new chapters and for providing invaluable feedback and suggestions that improved their content. Thanks also to Nick Horton, Stu Lipsitz, and Caitlin Ravichandran for their helpful suggestions and insightful comments on several chapters. Finally, we thank Steve Quigley and Susanne Steitz-Filler of Wiley, for their advice and encouragement during all stages of this project.

GARRETT M. FITZMAURICE  
NAN M. LAIRD  
JAMES H. WARE

Boston, Massachusetts  
May, 2011

## *Preface to First Edition*

Our goal in writing this book is to provide a rigorous and systematic description of modern methods for analyzing data from longitudinal studies. In recent years there have been remarkable developments in methods for longitudinal analysis. Despite these important advances, the methods have been somewhat slow to move into the mainstream. *Applied Longitudinal Analysis* bridges the gap between theory and application by presenting a comprehensive account of these methods in a way that is accessible to a wide range of readers.

The impetus for this book arose from teaching a graduate-level course on “Applied Longitudinal Analysis” at the Harvard School of Public Health. As course instructors, we were frustrated by the lack of a suitable textbook that adequately covered modern statistical methods for longitudinal analysis at a level accessible to a broad audience of researchers and graduate students in the health and medical sciences. We envision this book as a textbook for such a course and, subsequently, as a reference resource for researchers and graduate students. It is also suitable for graduate students in statistics and for statisticians already working in the health sciences, governmental health-related agencies, and the pharmaceutical industry. It is intended to allow a diverse group of statisticians, researchers, and graduate students in substantive fields to master modern methods for longitudinal data analysis.

The scope of this book is broad, covering methods for the analysis of diverse types of longitudinal data arising in the health sciences. The methods are presented in the setting of numerous applications to real data sets. Our main emphasis is on the practical rather than the theoretical aspects of longitudinal analysis. Twenty-five real data sets, drawn from studies in health-related fields, are used throughout the text and problem sets to illustrate the applications of longitudinal methods. These data sets can be downloaded from the website for the book: [www.biostat.harvard.edu/~fitzmaur/ala](http://www.biostat.harvard.edu/~fitzmaur/ala). Although the methods are applied to data sets drawn from the health sciences, they apply equally to other areas of application, for example, education, psychology, and other branches of the behavioral and social sciences.

Because longitudinal data are a special case of clustered data, albeit with a natural ordering of the measurements within a cluster, we include also a description of modern methods for analyzing clustered data, more broadly defined. Indeed, one of our goals is to demonstrate that methods for longitudinal analysis are, more or less, special cases of more general regression methods for clustered data. As a result a comprehensive understanding of longitudinal data analysis provides the basis for a broader understanding of methods for analyzing the wide range of clustered data that commonly arises in studies in the biomedical and health sciences.

The prerequisites for a course based on this book are an introductory course in statistics and a strong background in regression analysis. Some previous exposure to generalized linear models (e.g., logistic regression) would be helpful, although these models are reviewed in the text. An understanding of matrix algebra or calculus is not assumed; the reader will be gently introduced to only those aspects of vector and matrix notation necessary for understanding the matrix representation of regression models for longitudinal data. Because vectors and matrices are used to simplify notation, the reader is required to attain some basic facility with the addition and multiplication of vectors and matrices. Although we do not assume a high level of mathematical preparation, a willingness to read and consider mathematical ideas is required. More technical or mathematical sections of the book are marked with asterisks and may be omitted at first reading without loss of continuity.

To use the methods described in this book, appropriate statistical software is required. In general, the methods available via commercially available software lag behind the recent advances in statistical methods; longitudinal data analysis is not exceptional in this regard. Recently the introduction of new programs for analyzing multivariate and longitudinal data has made these methods far more accessible to practitioners and students. We use *SAS*, which is widely available, to

perform the analyses presented throughout the text. Illustrative *SAS* commands are included at the end of many of the chapters, with basic descriptions of their usage. Programming statements and computer output for the examples, prepared using *SAS*, can be downloaded from the website: [www.biostat.harvard.edu/~fitzmaur/ala](http://www.biostat.harvard.edu/~fitzmaur/ala). We selected *SAS* because all of the analyses we discuss can be performed using its procedures. Many of the methods can be carried out using alternative software packages (e.g., *S-Plus* and *Stata*) or special purpose programs (e.g., *BMDP5-V*) and this book can be supplemented with any one of them. Readers are encouraged to perform and verify the results of analyses using software of their choice. Because statistical software is constantly evolving, we anticipate that all of the methods we discuss will soon be available within most of the major statistical packages.

Throughout the text references have been kept to an absolute minimum. Instead, at the end of each chapter we include suggestions for further readings that provide more in-depth coverage of certain topics. We also include “bibliographic notes” that highlight key references in the mainstream statistical literature. Although many of our readers may find the latter references to be too technical, they are included to give due credit to those who have contributed to the statistical methods described in each chapter.

Finally, we would like to thank the many friends and colleagues who have helped us to write this book. A special word of thanks to Misha Salganik, for preparation of the diagrams and many helpful suggestions for improvement of graphical displays. We are especially grateful to Joe Hogan and Russell Localio, for reading a first draft and providing invaluable feedback, comments, and suggestions that improved the book. We would also like to thank Rino Bellocchio, Brent Coull, Nick Horton, Sharon-Lise Normand, Misha Salganik, Judy Singer, S. V. Subramanian, and Florin Vaida, for their insightful comments on several chapters. We are grateful to the students who have participated in the course on “Applied Longitudinal Analysis” at the Harvard School of Public Health since its inception; they have provided the impetus and motivation for writing this book. We gratefully acknowledge support from grant GM 29745 from the National Institutes of Health. The first author gratefully acknowledges support from the Junior Faculty Sabbatical Program at the Harvard School of Public Health; the support provided by a sabbatical created a unique opportunity to begin writing this book. Last, but not least, we thank Steve Quigley and Susanne Steitz of Wiley, for their advice and encouragement during all stages of this project.

GARRETT M. FITZMAURICE  
NAN M. LAIRD  
JAMES H. WARE

*Boston, Massachusetts*  
*March, 2004*

## *Acknowledgments*

Throughout this book we have used data sets drawn from published studies in health-related fields to exemplify important concepts in the analysis of longitudinal and clustered data. We are grateful to the following investigators for sharing their data with us: Graham Bentham, Doug Dockery, Brian Flay, Robert Greenberg, Keith Henry, Aviva Must, Elena Naumova, George Rhoads, Jan Schouten, Linda Van Matter, and Gwen Zahner.

We also thank the following publishers for permission to reproduce published data sets in print and electronic format: The American Statistical Association, Blackwell Publishing, Brooks/Cole (a division of Thomson Learning), CRC Press, Elsevier, Iowa State Press, Oxford University Press, and SAS Institute, Inc.

Finally, in all data sets used throughout this book, the original subject identification (ID) numbers have been deleted and replaced with new subject ID numbers, to ensure that the data sets cannot be linked to the original records.

## *Part I*

# *Introduction to Longitudinal and Clustered Data*

# *Chapter 1*

## *Longitudinal and Clustered Data*

### **1.1 INTRODUCTION**

Research on statistical methods for the design and analysis of human investigations expanded explosively in the second half of the twentieth century. Beginning in the early 1950s, the U.S. government shifted a substantial part of its research support from military to biomedical research. The legislative foundation for the modern National Institutes of Health (NIH), the Public Health Service Act, was passed in 1944 and NIH grew rapidly throughout the 1950s and 1960s. During these “golden years” of NIH expansion, the entire NIH budget grew from \$8 million in 1947 to more than \$1 billion in 1966. The NIH sponsored many of the important epidemiologic studies and clinical trials of that period, including the influential Framingham Heart Study (Dawber et al., 1951; Dawber, 1980).

The typical focus of these early studies was morbidity and, especially, mortality. Investigators sought to identify the causes of early death and to evaluate the effectiveness of treatments for delaying death and morbidity. In the Framingham Heart Study, participants were seen at two-year intervals. Survival outcomes during successive two-year periods were treated as independent events and modeled using multiple logistic regression. The successful use of multiple logistic regression in this setting, and the recognition that it could be applied to case-control data, led to widespread use of this methodology beginning in the 1960s. The analysis of time-to-event data was revolutionized by the seminal 1972 paper of D. R. Cox, describing the proportional hazards model (Cox, 1972). This paper was followed by a rich and important body of work that established the conceptual basis and the computational tools for modern survival analysis.

Although the design of the Framingham Heart Study and other cohort studies called for periodic measurement of the patient characteristics thought to be determinants of chronic disease, interest in the levels and patterns of change of those characteristics over time was initially limited. As the research advanced, however, investigators began to ask questions about the behavior of these risk factors. In the Framingham Heart Study, for example, investigators began to ask whether blood pressure levels in childhood were predictive of hypertension in adult life. In the Coronary Artery Risk Development in Young Adults (CARDIA) Study, investigators sought to identify the determinants of the transition from normotensive or normocholesterolemic status in early adult life to hypertension and hypercholesterolemia in middle age (Friedman et al., 1988). In the treatment of arthritis, asthma, and other diseases that are not typically life-threatening, investigators began to study the effects of treatments on the level and change over time in measures of severity of disease. Similar questions were being posed in every disease setting. Investigators began to follow populations of all ages over time, both in observational studies and clinical trials, to understand the development and persistence of disease and to identify factors that alter the course of disease development.

This interest in the temporal patterns of change in human characteristics came at a period when advances in computing power made new and more computationally intensive approaches to statistical analysis available at the desktop. Thus, in the early 1980s, Laird and Ware proposed the use of the EM algorithm to fit a class of linear mixed effects models appropriate for the analysis of repeated measurements (Laird and Ware, 1982); Jennrich and Schluchter (1986) proposed a variety of alternative algorithms, including Fisher-scoring and Newton–Raphson algorithms. Later in the decade, Liang and Zeger introduced the generalized estimating equations in the biostatistical literature and proposed a family of generalized linear models for fitting repeated observations of

binary and counted data (Liang and Zeger, 1986; Zeger and Liang, 1986). Many other investigators writing in the biomedical, educational, and psychometric literature contributed to the rapid development of methodology for the analysis of these “longitudinal” data. The past 30 years have seen considerable progress in the development of statistical methods for the analysis of longitudinal data. Despite these important advances, methods for the analysis of longitudinal data have been somewhat slow to move into the mainstream. This book bridges the gap between theory and application by presenting a comprehensive description of methods for the analysis of longitudinal data accessible to a broad range of readers.

## 1.2 LONGITUDINAL AND CLUSTERED DATA

The defining feature of longitudinal studies is that measurements of the same individuals are taken repeatedly through time, thereby allowing the direct study of change over time. The primary goal of a longitudinal study is to characterize the change in response over time and the factors that influence change. With repeated measures on individuals, one can capture within-individual change. Indeed, the assessment of within-subject changes in the response over time can only be achieved within a longitudinal study design. For example, in a cross-sectional study, where the response is measured at a single occasion, one can only obtain estimates of between-individual differences in the response. That is, a cross-sectional study may allow comparisons among sub-populations that happen to differ in age, but it does not provide any information about how individuals change during the corresponding period.

To highlight this important distinction between cross-sectional and longitudinal study designs, consider the following simple example. Body fatness in girls is thought to increase just before or around menarche, leveling off approximately 4 years after menarche. Suppose that investigators are interested in determining the increase in body fatness in girls after menarche. In a cross-sectional study design, investigators might obtain measurements of percent body fat on two separate groups of girls: a group of 10-year-old girls (a pre-menarcheal cohort) and a group of 15-year-old girls (a post-menarcheal cohort). In this cross-sectional study design, direct comparison of the average percent body fat in the two groups of girls can be made using a two-sample (unpaired)  $t$ -test. This comparison does not provide an estimate of the change in body fatness as girls age from 10 to 15 years. The effect of growth or aging, an inherently within-individual effect, simply cannot be estimated from a cross-sectional study that does not obtain measures of how individuals change with time. In a cross-sectional study the effect of aging is potentially confounded with possible cohort effects. Put in a slightly different way, there are many characteristics that differentiate girls in these two different age groups that could distort the relationship between age and body fatness. On the other hand, a longitudinal study that measures a single cohort of girls at both ages 10 and 15 can provide a valid estimate of the change in body fatness as girls age. In the longitudinal study the analysis is based on a paired  $t$ -test, using the difference or change in percent body fat within each girl as the outcome variable. This within-individual comparison provides a valid estimate of the change in body fatness as girls age from 10 to 15 years. Moreover, since each girl acts as her own control, changes in percent body fat throughout the duration of the study are estimated free of any between-individual variation in body fatness.

A distinctive feature of longitudinal data is that they are *clustered*. In longitudinal studies the clusters are composed of the repeated measurements obtained from a single individual at different occasions. Observations within a cluster will typically exhibit positive correlation, and this correlation must be accounted for in the analysis. Longitudinal data also have a temporal order; the first measurement within a cluster necessarily comes before the second measurement, and so on. The ordering of the repeated measures has important implications for analysis. There are, however, many studies in the health sciences that are not longitudinal in this sense but which give rise to data that are clustered or cluster-correlated. For example, clustered data commonly arise when intact groups are randomized to health interventions or when naturally occurring groups in the population are randomly sampled. An example of the former is group-randomized trials. In a group-randomized trial, also known as a cluster-randomized trial, groups of individuals, rather than each individual alone, are randomized to different treatments or health interventions. Data on the health outcomes of interest are obtained on all individuals within a group. Alternatively, clustered data can arise from random sampling of naturally occurring groups in the population. Families, households, hospital wards, medical practices, neighborhoods, and schools are all instances of naturally occurring clusters in the population that might be the primary sampling units in a study. Finally, clustered data can arise when data on the health outcome of interest are simultaneously obtained either from multiple raters or from

different measurement instruments.

In all these examples of clustered data, we might reasonably expect that measurements on units within a cluster are more similar than the measurements on units in different clusters. The degree of clustering can be expressed in terms of correlation among the measurements on units within the same cluster. This correlation invalidates the crucial assumption of independence that is the cornerstone of so many standard statistical techniques. Instead, statistical models for clustered data must explicitly describe and account for this correlation. Because longitudinal data are a special case of clustered data, albeit with a natural ordering of the measurements within a cluster, this book includes a description of modern methods of analysis for clustered data, more broadly defined. Indeed, one of the goals of this book is to demonstrate that methods for the analysis of longitudinal data are, more or less, special cases of more general regression methods for clustered data. As a result a comprehensive understanding of methods for the analysis of longitudinal data provides the basis for a broader understanding of methods for analyzing the wide range of clustered data that commonly arises in studies in the biomedical and health sciences.

The examples described above consider only a single level of clustering, for example, repeated measurements on individuals. More recently investigators have developed methodology for the analysis of multilevel data, in which observations may be clustered at more than one level. For example, the data may consist of repeated measurements on patients clustered by clinic. Alternatively, the data may consist of observations on children nested within classrooms, nested within schools. Although the analysis of multilevel data is not the primary focus of this book, multilevel data are discussed in Chapter 22.

Interest in the analysis of longitudinal and multilevel data continues to grow. New and more flexible models have been developed and advances in computation, such as Markov chain Monte Carlo (MCMC) methods, have allowed greater flexibility in model specification. Moreover, improvements in statistical software packages, especially SAS, Stata, SPSS, R, and S-Plus, have made these models much more accessible for use in routine data analysis. Despite these advances, however, methods for the analysis of longitudinal data are not widely used and are seen to be accessible only to statisticians with specialized expertise.

We believe that the methodology for the analysis of longitudinal data can be much more widely understood and applied. It is our hope that this book will help make that possible. It provides a comprehensive introduction to methods for the analysis of longitudinal data, written for a reader with a basic knowledge of statistics and a strong background in regression analysis. The book does not require a high level of mathematical preparation but does assume a willingness to read and consider mathematical ideas.

## **1.3 EXAMPLES**

To highlight some of the distinctive features of longitudinal and clustered data, we introduce four examples drawn from studies in the biomedical sciences. These four examples will be used later in the book to illustrate different analytic approaches. Additional examples, also drawn from studies in the biomedical and health sciences, will be introduced in later chapters of the book.

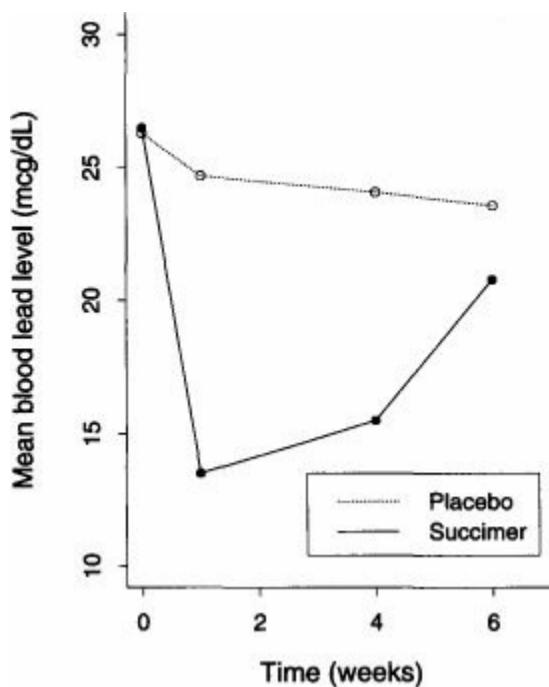
## 1.3.1 Treatment of Lead-Exposed Children (TLC) Trial

Exposure to lead can produce cognitive impairment, especially among young children and infants. A young child exposed to high levels of lead may experience various adverse health effects, including hyperactivity, hearing or memory loss, learning disabilities, and damage to the brain and nervous system. Although the use of lead as an additive in gasoline has been discontinued, at least in the United States, resulting in a dramatic reduction in airborne lead levels, a small percentage of children continue to be exposed to lead at levels that can produce impairment. Much of this exposure is due to deteriorating lead-based paint (e.g., chipping and peeling paint) in older homes. Lead was used as a pigment and drying agent in “alkyd” oil-based paint. While the United States government banned the use of lead-based paint in housing in 1978, many homes built before 1978 contain lead-based paint. When lead-based paint deteriorates, it becomes lead paint chips, which can be eaten by young children, and lead-contaminated paint dust, which can be ingested by young children during normal teething and hand-to-mouth behavior. The U.S. Centers for Disease Control and Prevention (CDC) has concluded that children with blood lead levels above 10 micrograms per deciliter ( $\mu\text{g}/\text{dL}$ ) of whole blood are at risk of adverse health effects.

Lead poisoning in children is treatable in the sense that there are medical interventions, known as chelation treatments, that can help a child to excrete the lead that has been ingested. Until recently chelation treatment of children with high levels of blood lead was administered by injection and required hospitalization. A new chelating agent, succimer, enhances urinary excretion of lead and has the distinct advantage that it can be given orally, rather than by injection. In the 1990s the *Treatment of Lead-Exposed Children (TLC) Trial Group* conducted a placebo-controlled, randomized trial of succimer in children with confirmed blood lead levels of 20 to 44  $\mu\text{g}/\text{dL}$ , levels well above the CDC’s threshold for concern about the adverse health effects of exposure to lead (Treatment of Lead-Exposed Children (TLC) Trial Group, 2000; Rogan et al., 2001). The children were aged 12 to 33 months at enrollment and lived in deteriorating inner city housing. The mean age of the children at randomization was 2 years and the mean blood lead level was 26  $\mu\text{g}/\text{dL}$ . Children received up to three 26-day courses of succimer or placebo and were followed for 3 years.

[Table 1.1](#) presents data on blood lead levels at baseline, week 1, week 4, and week 6 for 10 randomly selected children from the study. The mean blood lead levels at each measurement occasion for a random subset of 100 children, broken down by treatment group, are presented in [Table 1.2](#). As expected, due to randomization, the mean response at baseline is similar in the two treatment groups. However, there are discernible differences in the patterns of change in the mean response over time. A graphical presentation of the mean blood lead levels at each occasion is displayed in [Figure 1.1](#). Note that at week 1 there appears to be a dramatic drop in initial blood lead levels among the children treated with succimer. However, this is followed by a rebound in blood lead levels, as lead stored in the bones and tissues is mobilized and a new equilibrium is achieved. In contrast, for the children treated with placebo, the trend in the mean response over time is relatively flat.

[Fig. 1.1](#) Plot of mean blood lead levels at baseline, week 1, week 4, and week 6 in the succimer and placebo groups.



**Table 1.1** Blood lead levels ( $\mu\text{g}/\text{dL}$ ) at baseline, week 1, week 4, and week 6 for 10 randomly selected children from the TLC trial.

ID	Group <sup>a</sup>	Baseline	Week 1	Week 4	Week 6
79	P	30.8	26.9	25.8	23.8
8	S	26.5	14.8	19.5	21.0
44	S	25.8	23.0	19.1	23.2
11	P	24.7	24.5	22.0	22.5
69	S	20.4	2.8	3.2	9.4
29	S	20.4	5.4	4.5	11.9
46	P	28.6	20.8	19.2	18.4
13	P	33.7	31.6	28.5	25.1
74	P	19.7	14.9	15.3	14.7
53	P	31.1	31.2	29.2	30.1

<sup>a</sup> P = placebo; S = succimer.

**Table 1.2** Mean blood lead levels (and standard deviation) at baseline, week 1, week 4, and week 6 for children from the TLC trial.

Group	Baseline	Week 1	Week 4	Week 6
Succimer	26.5 (5.0)	13.5 (7.7)	15.5 (7.8)	20.8 (9.2)
Placebo	26.3 (5.0)	24.7 (5.5)	24.1 (5.8)	23.6 (5.6)

## 1.3.2 Muscatine Coronary Risk Factor Study

In 1998 the American Heart Association (AHA) announced that obesity had been added to the AHA's list of major preventable risk factors for coronary heart disease. These major preventable risk factors include smoking, high blood cholesterol, high blood pressure, and sedentary lifestyle. Unlike risk factors that cannot be altered, such as heredity, increasing age, and being male, obesity is a risk factor that many individuals can alter and control. The medical definition of obesity is quite simple: an excess of body fat. Obesity is primarily caused by consuming too many calories and not getting enough physical exercise. Obesity can lead to higher blood cholesterol and triglyceride levels, lower HDL cholesterol (HDL cholesterol, the "good" cholesterol, has been linked to lower risk of coronary heart disease), and higher blood pressure. Thus obesity can contribute to higher coronary risk in a variety of different ways.

Public health scientists now accept that obesity is a chronic disease, just like high blood pressure or high blood cholesterol. Its causes are a complex, individualized combination of genetics, behavior, and lifestyle. There is also increased awareness that obese children are at increased risk for obesity as adults.

In 1970 researchers from the University of Iowa began to examine the links between child and adult coronary health. Of particular interest were the associations between coronary risk factors in youth and coronary disease in adults. The Muscatine Coronary Risk Factor (MCRF) study, a longitudinal survey of school-age children in Muscatine, Iowa, had the goal of examining the development and persistence of risk factors for coronary disease in children (Woolson and Clarke, 1984; Lauer et al., 1997). In the MCRF study, weight and height measurements of five cohorts of children, initially aged 5–7, 7–9, 9–11, 11–13, and 13–15 years, were obtained biennially from 1977 to 1981. Data were collected on 4856 boys and girls. On the basis of a comparison of their weight to age-gender specific norms, children were classified as obese or not obese. One objective was to determine whether the prevalence of obesity increases with age and whether patterns of change in obesity are the same for boys and girls.

A summary of the obesity data for children in one of the five cohorts, who were 7–9 years old in 1977, is presented in [Table 1.3](#). Because all the variables are discrete, the data can be summarized as counts in a contingency table. For example, the first 8 rows of [Table 1.3](#) provide a count of the number of children with each of the 8 (or  $2^3$ ) possible sequences of binary responses over the three measurement occasions. A similar table could be constructed for each of the remaining four cohorts of children. Note that although each child was eligible to participate in all three surveys, the data are incomplete for many children. Less than 40% of the children provided complete data at all three measurement occasions. For convenience, in [Table 1.3](#) the missingness of obesity is treated as a third category of the obesity status variable.

**Table 1.3** Obesity status of cohort of children, aged 7–9 at entry, from the Muscatine study.

Gender	Child's Obesity Status <sup>a</sup>			Count
	1977	1979	1981	
<b>Males</b>				
None missing	1 1 1 1 0 0 0 0	1 1 0 0 1 1 0 0	1 0 1 0 1 0 1 0	20 7 9 8 8 8 15 150
Missing time 1	*	1 1 0 *	1 0 1 0	13 3 2 42
Missing time 2	1 1 0 0	*	1 0 1 0	3 1 6 16
Missing time 3	1 1 0 0	1 0 1 0	*	11 1 3 38
Missing times 1, 2	*	*	1	14
Missing times 1, 3	*	*	0	55
Missing times 2, 3	1 0	*	*	33 7 45
<b>Females</b>				
None missing	1 1 1 1 0 0 0 0	1 1 0 0 1 1 0 0	1 0 1 0 1 0 1 0	21 6 6 2 19 13 14 154
Missing time 1	*	1 1 0 *	1 0 1 0	8 1 4 47
Missing time 2	1 1 0 0	*	1 0 0 1	4 0 16 3
Missing time 3	1 1 0 0	1 0 1 0	*	11 1 3 25
Missing times 1, 2	*	*	1	13
Missing times 1, 3	*	*	0	39
Missing times 2, 3	1 0	*	*	5 23 7 47

<sup>a</sup> 1 = Obese; 0 = Not Obese; \* = Missing.

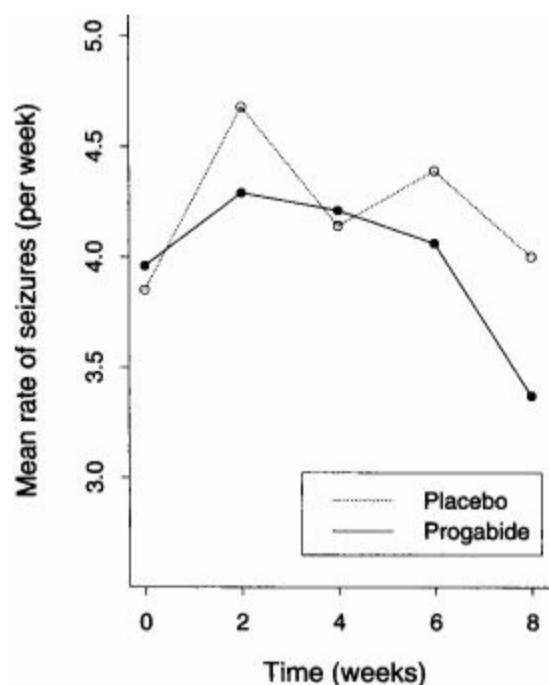
### 1.3.3 Clinical Trial of an Anti-epileptic Drug

Epilepsy is a chronic neurologic disorder that may result from brain injury, developmental malformation, or a genetic abnormality. It is characterized by recurrent seizures caused by sudden, excessive electrical activity in the brain. Seizures are classified as generalized, in which the electrical discharge occurs throughout the brain, and partial onset, wherein the electrical activity is localized.

Data for the third example come from a placebo-controlled clinical trial of 59 epileptics conducted by Leppik et al. (1987). Patients with partial seizures were enrolled in a randomized clinical trial of the anti-epileptic drug, progabide. Participants in the study were randomized to either progabide or a placebo, as an adjuvant to the standard anti-epileptic chemotherapy. Progabide is an anti-epileptic drug whose primary mechanism of action is to enhance gamma-aminobutyric acid (GABA) content; GABA is the primary inhibitory neurotransmitter in the brain.

Prior to receiving treatment, baseline data on the number of epileptic seizures during the preceding 8-week interval were recorded. Counts of epileptic seizures during 2-week intervals before each of four successive post-randomization clinic visits were recorded. The average rates of seizures (per week) at baseline and in the four post-randomization visits are presented in [Table 1.4](#). A graphical presentation of the average rates of seizures at each occasion in the progabide and placebo groups is displayed in [Figure 1.2](#). The main goal of the study was to compare the changes in the average rates of seizures in the two groups.

**Fig. 1.2** Mean rate of seizures (per week) at baseline, week 2, week 4, week 6, and week 8 in the progabide and placebo groups.



**Table 1.4** Mean rate of seizures per week (and standard deviation) at baseline, week 2, week 4, week 6, and week 8 in the clinical trial of progabide.

Group	Baseline	Week 2	Week 4	Week 6	Week 8
Progabide	3.96 (3.5)	4.29 (9.1)	4.21 (5.9)	4.06 (7.0)	3.37 (5.6)
Placebo	3.85 (3.3)	4.68 (5.1)	4.14 (4.1)	4.39 (7.3)	4.00 (3.8)

### 1.3.4 Connecticut Child Surveys

There is now accumulating evidence that the rates of psychiatric disorders in children are substantial, with reported population prevalence rates of childhood psychopathology ranging from 12% to 22%. However, children are considered to be unreliable in reporting on their own psychopathology. As a result many contemporary surveys of childhood psychopathology use proxy informants, usually a child's parent (or primary caregiver) and teacher, to report on the child's psychiatric status. In numerous studies the agreement among multiple informant reports on the child's psychopathology has been found to be poor. It is thought that much of this disagreement is less a result of the unreliability of the informant reports than of true differences in children's behaviors and emotions across different situations and settings, notably in the home and school. A central issue in studies of risk factors for childhood psychopathology is utilization of the information obtained about the child's mental health status from multiple sources or informants.

Data for our example come from two parallel epidemiological surveys that assessed the mental health and service needs of children, aged 6 to 11, in rural and urban communities in Connecticut (Zahner et al., 1992, 1993). The first survey, the New Haven Child Survey (NHCS), was conducted in 1986 and 1987 in New Haven, Connecticut, a predominantly minority metropolitan center. The second survey, the Eastern Connecticut Child Survey (ECCS) was conducted in 1988 and 1989 and replicated the NHCS in a non-metropolitan planning region covering the eastern third of Connecticut. The two studies used comparable survey procedures. In particular, they used parallel questionnaires designed to be self-administered by the children's parents and teachers. Children's emotional and behavioral problems were assessed with the Child Behavior Checklist (CBCL) and the Teacher's Report Form (TRF), 118-item symptom inventories covering problems commonly seen in child guidance clinics. The CBCL and TRF scales do not provide diagnoses of psychiatric disorders; instead, they provide broad-band measures of emotional (or "internalizing") and behavioral (or "externalizing") disturbance. The CBCL and TRF scale scores can be dichotomized at published clinical cut-points.

Thus the New Haven Child Survey and the Eastern Connecticut Child Survey provided both a parent's and a teacher's report of psychiatric disturbance in the child as assessed by parallel forms of a standardized psychiatric symptom checklist. These data provide multiple source (here, from two sources: the parent and teacher) information on the psychiatric outcome variable of interest. Of note, these data are cross-sectional but the two sources of information about each child's psychopathology are likely to be positively correlated. Thus data from the Connecticut Child Surveys are an example of clustered, but not longitudinal, data. In this setting, unlike a typical longitudinal study, the major interest of the analysis is not in changes in the response over time. Instead, the major focus of the analysis is on the effects of subject-specific covariates on the outcome.

[Table 1.5](#) displays social and demographic characteristics of the children and the overall rates of externalizing disturbance as determined by CBCL and TRF scale scores in the clinical range.

**Table 1.5** Frequency distribution for variables from the Connecticut Child Surveys.

Variables	Count	Percent
<i>Parent informant (N = 2501)</i>		
Externalizing		
0 = Normal	2112	84
1 = Borderline/clinical	389	16
<i>Teacher informant (N = 1428)</i>		
Externalizing		
0 = Normal	1159	81
1 = Borderline/clinical	269	19
Area		
1 = Rural	874	35
2 = Suburban	428	17

3 = Small city	386	15
4 = Large city	813	33
Single parent		
0 = No	1982	79
1 = Yes	519	21
Child's health		
0 = Good health	1329	53
1 = Fair/bad health	1172	47
Child's gender		
0 = Female	1294	52
1 = Male	1207	48

The four examples considered in this section differ in terms of outcome variable, study design, and goals or objectives of the analysis. In the first example from the TLC trial, the outcome variable, blood lead level, is continuous. In the second example from the MCRF study, the outcome variable, obesity status, is binary. In the third example from the clinical trial of progabide, the outcome variable is a count. These three examples illustrate the diverse types of longitudinal data that arise in the health and medical sciences. A notable feature of the second example is the amount of missing data. Missing data are a common problem in longitudinal studies in the health sciences. As we will discuss in later chapters, one will need to examine the reasons for any missingness to determine the validity of inferences about changes in the response over time. Next, consider the design of these studies. The first and third examples are experiments, where the treatments have been chosen by the investigators and randomly assigned to the study participants. The second example is an observational study where the study participants are followed forward in time to observe the outcome variable at future time points; however, unlike the randomized clinical trial, the investigators cannot directly control the comparability of groups (here, males and females). While the first three examples involve longitudinal study designs, the fourth example is a cross-sectional observational study. In the Connecticut Child Surveys, variables are measured at a single time point on a sample of children. Because information on the outcome variable of interest is obtained from two sources (the parent and teacher), these data are also clustered. Finally, we note that the goals of the analysis are similar for the first three examples: characterize the change in the outcome variable over time and the factors that influence change. In the fourth example, however, the objective of the analysis is not to characterize change in the outcome variable over time. Instead, the goal is to examine the effects of subject-specific covariates on the outcome. In later chapters we describe modern methods for analyzing diverse types of longitudinal data arising from both experiments and observational studies. Because longitudinal data are a special case of clustered data, we also describe methods of analysis for clustered data, more broadly defined.

# 1.4 REGRESSION MODELS FOR CORRELATED RESPONSES

In the last 30 years we have seen remarkable advances in methods for analyzing longitudinal and clustered data. In particular, we now have a broad and flexible class of models for correlated data based on a regression paradigm. Indeed, all the methods that are described in later chapters can be thought of as regression models for correlated responses. In this section we provide motivation for the regression paradigm for correlated responses.

Regression models are widely used and provide a very general and versatile approach for analyzing data. Our use of the term “regression model” here is not strictly limited to the standard linear regression model for a continuous response variable. Instead, we use this term more broadly to refer to any model that describes the dependence of the mean of a response variable on a set of covariates in terms of some form of regression equation. While the simplest case is the familiar linear regression model for a continuous response variable, there are many possible generalizations. For example, regression models have been developed for other response variables, such as binary responses or counts. For the binary response variable, linear logistic regression has been widely used for many applications. For counts, Poisson or log-linear regression is often appropriate. Another important generalization is to observations that cannot be assumed to be statistically independent of one another, that is, regression models for correlated responses. In later chapters we consider both kinds of generalizations of the standard linear regression model.

Note that the term “linear” has appeared in all three of the examples of regression models considered so far. Linearity in this setting has a very precise meaning and refers to the fact that all of these models for the mean (or some transformation of the mean) are linear in the regression parameters. For example, letting  $Y$  denote the response variable and  $X$  a covariate, the following three models for the mean response

$$E(Y|X) = \beta_1 + \beta_2 X,$$

$$E(Y|X) = \beta_1 + \beta_2 \log(X),$$

and

$$E(Y|X) = \beta_1 + \beta_2 X + \beta_3 X^2,$$

are all cases where the mean is linear in the regression parameters (where  $E(Y|X)$  denotes the *conditional* mean or expectation of  $Y$  given  $X$ ). All three models are linear in the regression parameters, even if the latter two are non-linear in the covariate. In this book we only consider models where the mean response, or some suitable transformation of the mean response (e.g., log transformation in Poisson regression), is linear in the regression parameters. We do not consider models that are fundamentally non-linear in the regression parameters. For example, the following two models

$$E(Y|X) = \beta_1 + e^{\beta_2 X},$$

and

$$E(Y|X) = \frac{\beta_1}{1 + \beta_2 e^{-\beta_3 X}},$$

are cases where the mean is non-linear in the regression parameters. However, we remind the reader that our focus on models that are linear in the regression parameters does not preclude relationships between the mean response and covariates that are curvilinear or non-linear. This type of non-linearity can be accommodated by taking appropriate transformations of the mean response (e.g., log transformation in Poisson regression) and the covariates (e.g.,  $\log(\text{dose})$ ), and/or by including polynomials. For example, a quadratic trend in the mean response over time can be incorporated by including both time and time<sup>2</sup> in the regression model. The inclusion of transformed covariates in no way violates the “linearity” of the regression model; that is, the model is still linear in the regression parameters.

As noted earlier, we use the term “regression model” to refer to any model that describes the

dependence of the response variable on a set of covariates in some form of regression equation. In particular, the regression parameters express how the mean of the response variable depends on the covariates. For example, in the case of the linear regression model for a continuous response, the regression coefficients express the dependence of the mean of the outcome in terms of a linear combination of the covariates. In the linear logistic model for a binary response, the regression coefficients express the dependence of the log odds of a positive response in terms of a linear combination of the covariates. Note, however, that the log odds is simply a non-linear transformation of the mean or probability of a positive response. Thus in both cases the mean of the response variable, or some appropriate transformation of the mean, is related to a linear combination of the covariates.

One appealing aspect of the regression paradigm concerns the nature of the explanatory variables. A feature of the regression modeling approach is that it can incorporate mixtures of discrete and continuous covariates in a relatively seamless fashion. That is, the covariates can be continuous (and often referred to as quantitative), such as body weight, age, time, and dose. Furthermore the mean response, or any suitable transformation of the mean, can be related to a continuous covariate in a curvilinear or non-linear fashion by simply taking an appropriate transformation of the covariate or by the inclusion of polynomials (e.g., time and time<sup>2</sup>). Alternatively, the covariates can be discrete (or qualitative), such as gender and treatment group. Finally, regression models can include mixtures of discrete and continuous covariates, and products among them. As a result, within a regression paradigm, it is no more difficult to analyze longitudinal data arising from a carefully designed experiment with a single qualitative covariate or factor (e.g., a randomized placebo-controlled longitudinal clinical trial) than from an observational study where there are many covariates, some of which are discrete, the others continuous. Of note, in the latter case, regression models can often be used to distinguish within- and between-subject trends in the response (e.g., “longitudinal” versus “cross-sectional” effects of age); this topic will be discussed in greater depth in later chapters.

Regression models can usually be formulated in such a way that certain regression parameters have interpretations that bear directly on the scientific question of main interest. For example, in a regression model for data from a longitudinal clinical trial, a particular regression coefficient can be given an interpretation in terms of the constant rate of change in the mean response over time in one of the treatment groups. Alternatively, the absence (or setting to zero) of a particular regression coefficient can be given an interpretation in terms of two treatment groups having the same underlying rate of change in the response variable over time.

So far we have emphasized that it is not necessary to distinguish whether the covariates are continuous or discrete (or a mixture of the two) within a regression paradigm. However, from a purely historical perspective, linear models for a continuous response with only discrete covariates have often been referred to as *analysis of variance* (ANOVA) models. In contrast, linear models for a continuous response with only continuous covariates have often been referred to as *linear regression* models. Indeed, some textbooks and courses in statistics present linear regression and analysis of variance as almost distinct analytic procedures. A large part of the reason for this arbitrary distinction is historical. Analysis of variance had its earliest roots in agricultural applications, especially carefully designed experiments where the responses (e.g., crop yield) could be indexed by one or more classifying factors (e.g., plot, crop variety) or qualitative experimental factors (e.g., different types of fertilizers). In contrast, linear regression was initially developed for the analysis of observational data. Some of the earliest applications of linear regression can be traced back to astronomy. By their very nature the data arising from studies in astronomy were purely observational (e.g., the positions and magnitudes of the heavenly bodies) and not the product of experimental manipulations. As a result of their somewhat different historical roots, ANOVA and linear regression have often been presented as almost distinct procedures, intended for the analysis of data arising from studies that differ in design (experimental versus observational) and the nature of the covariates (discrete versus continuous). Later it was recognized that linear regression is a very general model that incorporates analysis of variance as a special case.

Thus, although many of the commonly used statistical models for correlated data were originally

developed for data arising from studies that differed in design, aims, and the nature of the covariates, almost all of these developments fall within the regression paradigm for correlated data. So from a purely pedagogical perspective, it is not necessary to distinguish methods for analyzing longitudinal or correlated data arising from observational studies and from studies with experimental designs. From this point of view, we have purposely chosen not to focus on many of the early developments in methodology for analyzing correlated data, for example, the repeated measures ANOVA and multivariate analysis of variance (MANOVA). Instead, we focus on a more general and versatile regression paradigm that encompasses most, if not all, of the earlier developments as special cases but can also handle all of the complexities that arise in applications. When viewed as special cases within the regression paradigm, the underlying (and often unrealistic) assumptions made by many of the earliest methods for analyzing correlated data are more readily understood.

In summary, we view the regression paradigm as a very flexible and versatile approach for analyzing longitudinal and correlated data arising from many different types of studies. Regression models can provide a parsimonious description or explanation of how the mean response in a longitudinal study changes with time, and how these changes are related to covariates of interest. Thus our use of regression models is primarily intended for descriptive purposes, that is, for determining the most salient aspects of patterns of change in the mean response. While this does not necessarily preclude their use as a possible explanation of the underlying probabilistic data generating mechanism that might have produced the repeated responses, the latter is not considered to be the main focus of the analysis. Instead, our primary goal is to provide a simple description of the discernible patterns of change in the response over time, and their relation to covariates, via regression coefficients that bear directly on the scientific questions of main interest.

# 1.5 ORGANIZATION OF THE BOOK

The book is organized into five main parts. The first part, consisting of Chapters 1 and 2, provides the reader with an overview of the most salient aspects of longitudinal data. In Chapter 2, we introduce some notation and many of the analytic issues that arise with longitudinal data. We discuss the main features that distinguish longitudinal data from cross-sectional data. We highlight the major goals and objective of longitudinal analysis. We consider the aspect of longitudinal data that complicates their analysis, namely the correlation among repeated measures on the same individuals. We provide some intuition for how and why the correlation arises in longitudinal data and the potential consequences of ignoring it in the analysis.

The second part, consisting of Chapters 3 through 10, focuses on methods for analyzing longitudinal data when the response variable is continuous and assumed to have an approximate multivariate normal (or Gaussian) distribution. In Chapter 3, we introduce a general linear regression model for longitudinal data. We present a broad overview of different approaches for modeling the mean response over time and for accounting for the correlation among repeated measures on the same individual. These topics are discussed in much greater depth in subsequent chapters. In Chapter 4, we discuss estimation, via the method of maximum likelihood (ML), and inference concerning the regression coefficients and the covariance among the repeated measures. Longitudinal data present us with two aspects of the data that require modeling: the mean response over time and the covariance among repeated measures on the same individuals. In Chapters 5 and 6, the emphasis is on modeling the mean response. Two main approaches are distinguished: the analysis of response profiles (Chapter 5) and parametric or semiparametric curves (Chapter 6). In Chapter 7, we discuss models for the covariance in longitudinal data and develop an overall modeling strategy that takes account of the interdependence between the models for the mean and covariance. Chapter 8 introduces a very flexible class of models for analyzing longitudinal data known as linear mixed effects models. These models assume that some subset of the regression parameters vary randomly from one individual to another, thereby accounting for sources of natural heterogeneity in the population. Specifically, the mean response is modeled as a combination of fixed effects that are assumed to be shared by all individuals, and random effects that are unique or specific to a particular individual. In Chapter 9, we discuss an alternative, but closely related, class of regression models for longitudinal data known as linear “fixed effects” models. These models treat the subject-specific effects as fixed rather than random. We review the main features of linear fixed effects models for longitudinal data and discuss their potential advantages and disadvantages relative to linear mixed effects models. In Chapter 10, we discuss residual diagnostics for assessing the adequacy of models for longitudinal data and for detecting outlying observations and/or outlying individuals.

The chapters in the second part of the book cover many of the well-established methods for the analysis of longitudinal data and provide the foundation for future chapters that focus on discrete response variables (e.g., repeated binary responses and repeated count data). The third part, consisting of Chapters 11 through 16, focuses on methods for analyzing longitudinal data with outcomes that are not continuous. When the response is discrete, linear models are no longer appropriate for relating the mean to covariates. Instead, we consider extensions of generalized linear models for longitudinal data. In Chapter 11, we review the most salient features of generalized linear models for a single, univariate response; in later chapters, we discuss how generalized linear models can be extended to handle longitudinal responses. In generalized linear models a suitable non-linear transformation of the mean response is related to the covariates. However, this non-linearity raises some additional issues concerning the interpretation of the regression coefficients. In Chapters 12 through 15, we present two classes of models for analyzing discrete longitudinal data that account for the correlation among repeated measures in fundamentally different ways. In Chapter 16, we compare and contrast these two classes of models. One of the underlying themes emphasized in Chapters 12 through 16 concerns how different models for discrete longitudinal data have somewhat different targets of inferences. Thus, to ensure that the regression parameters bear directly on the question of scientific interest, greater care is needed in the choice of model for discrete longitudinal

data.

The fourth part of the book, consisting of Chapters 17 and 18, addresses the issue of missing data in longitudinal studies. In Chapter 17, we review the assumptions about missing data required to ensure that the methods discussed in earlier chapters provide valid inferences. Two methods for handling missing data, multiple imputation and inverse probability weighted methods, are discussed in detail in Chapter 18.

The final part of the book, consisting of Chapters 19 through 22, focuses on a number of advanced topics. In Chapter 19, we discuss smoothing methods for longitudinal analysis that allow greater flexibility for the form of the relationship between the mean response and the covariates. This chapter focuses on the connection between penalized splines and linear mixed effects models. Chapter 20 considers the design of a longitudinal study, focusing on the determination of sample size and power. In Chapter 21, we discuss regression models for repeated measures and related designs and emphasize how the methods discussed in earlier chapters can be applied in these settings. In Chapter 22, we present an overview of methods for analyzing multilevel data. Chapters 21 and 22 demonstrate how regression models for longitudinal data are special cases of general regression models for correlated data, more broadly defined.

## **1.6 FURTHER READING**

The presentation of methodology for the analysis of longitudinal data in subsequent chapters assumes that the reader has a basic knowledge of statistics and a strong background in regression analysis. A useful review of introductory statistical principles and methods, targeted at applied researchers, can be found in the books by Pagano and Gauvreau (2000) and Altman (1990). A comprehensive overview of regression concepts can be found in Kleinbaum et al. (1998) and Gelman and Hill (2007); a more advanced presentation of similar topics can be found in Neter et al. (1996).

# *Chapter 2*

## *Longitudinal Data: Basic Concepts*

### **2.1 INTRODUCTION**

In this chapter we present a broad overview of the main objectives of longitudinal analysis and some of the defining features of longitudinal data. Our primary goal is to emphasize that the major focus of the analysis of longitudinal data is on the assessment of within-individual changes in the response variable over time. That is, longitudinal analysis is concerned with estimating how individuals change throughout the duration of the study and examining the factors that influence heterogeneity among individuals in how they change over time. We also review the most salient features of longitudinal study designs, introduce some notation for longitudinal data, and highlight the main aspects of longitudinal data that complicate their analysis. Many of the concepts and issues introduced here will be discussed in much greater depth in later chapters of the book.

## 2.2 OBJECTIVES OF LONGITUDINAL ANALYSIS

In the health sciences, longitudinal studies play an important role in enhancing our understanding of the development and persistence of disease. There is much natural heterogeneity among individuals in terms of how diseases develop and progress. This heterogeneity is due to genetic, environmental, social, and behavioral factors. A longitudinal study design permits the discovery of individual characteristics that can explain these inter-individual differences in changes in health outcomes over time.

The distinguishing feature of longitudinal studies is that the study participants are measured repeatedly throughout the duration of the study, thereby permitting the direct assessment of changes in the response variable over time. In cross-sectional studies, where measurements are obtained at only a *single* point in time, it is not possible to assess individual changes on the basis of a single snapshot of the individual's response taken at a given time. Thus the defining feature of a longitudinal study is that two or more observations of the response variable, taken at different times, are made on at least some of the study participants. Typically, although not always, longitudinal study designs call for a fixed number of repeated measurements to be made on all study participants at a set of common time points. The occasions of measurement are not necessarily distributed evenly throughout the duration of the study.

By obtaining measurements of the same individuals repeatedly through time, longitudinal studies can address fundamental questions concerning the assessment of within-individual changes in the response variable. The main goal, indeed the *raison d'être*, of a longitudinal study, is to characterize the change in the response over time. While the measurement of within-individual changes is a fundamental objective of a longitudinal study, it is also of interest to determine whether these within-individual changes in the response are related to selected covariates. For example, in the *Treatment of Lead-Exposed Children Trial*, introduced in Chapter 1, repeated measures of blood lead levels were obtained at baseline (or week 0), week 1, week 4, and week 6, thereby allowing assessment of within-individual changes in blood lead levels over a six-week period. In this study it was not simply of interest to describe the overall pattern of within-individual changes in blood lead levels over time but also to relate these changes to the assigned treatment (placebo versus succimer).

In its most elementary form, a measure of the observed within-individual change in the response can be conceptualized in terms of simple "change scores" or "difference scores," for example, the differences between post-treatment and pre-treatment measurements of the response. The main objective of a longitudinal analysis is to describe trends in these within-individual changes in the response and to relate these changes to selected covariates (e.g., treatment group). This simple notion of within-individual change extends naturally from "difference scores" to more general "response trajectories" over time. For example, a "difference score" happens to be proportional to the slope (or constant rate of change) of a linear response trajectory. However, other kinds of response trajectories, for example, piecewise linear or curvilinear, can be used to parsimoniously smooth and summarize within-individual changes in the response throughout the duration of the study. In either case the fundamental ideas remain the same: we want to assess and describe within-individual changes in the response over time via comparison of measurements on the same individual taken later in time with those taken earlier.

A longitudinal analysis of within-individual changes proceeds in two conceptually distinct stages. First, within-individual change in the response is characterized in terms of some appropriate summary of the changes in the repeated measurements on each individual during the period of observation (e.g., using "difference scores" or some form of "response trajectory"). Second, these estimates of within-individual changes are then related to inter-individual differences in selected covariates. Although these two stages of the analysis are conceptually distinct, they can be combined in a statistical model for longitudinal data. That is, a single statistical model for longitudinal data can be used both to capture how individuals change over time and to relate within-individual changes in

the response to selected covariates.

For example, in the *Treatment of Lead-Exposed Children Trial* the investigators were interested in assessing changes in blood lead levels over time. In particular, they wanted to determine whether chelation treatment with succimer reduced blood lead levels over time relative to any changes in the placebo group. This study question can be addressed in an analysis that compares the two treatment groups in terms of the differences between post-treatment and pre-treatment measurements of blood lead levels. Although the major objective of the analysis is quite clear, there are many ways to construct and test hypotheses concerning treatment effects on changes in blood lead levels over time. For instance, the two treatment groups can be compared in terms of all post-treatment changes in the mean blood lead levels from baseline (or pre-treatment). Alternatively, the two treatment groups can be compared in terms of the rate of decline of blood lead levels over time, where the rate of decline is expressed in terms of a slope. Thus, although the scientific question of interest has a seemingly simple formulation in terms of whether changes in blood lead levels are affected by treatment, there are many different ways to proceed with a longitudinal analysis of these data. The choice of one analytic approach over another will usually depend on statistical considerations (e.g., issues of precision), the design of the study, and the specific scientific question of interest. These are topics that will be discussed in more detail in later chapters of the book.

Finally, it is an inescapable fact that the assessment of within-subject changes in the response over time can be achieved only within a longitudinal study design. A cross-sectional study simply cannot estimate how individuals change over time since the response is measured at only a single occasion. A longitudinal study can estimate how individuals change and also do so with great precision because each individual acts as his or her own control. By comparing each individual's responses at two or more occasions, a longitudinal analysis can remove extraneous, but unavoidable, sources of variability among individuals. The key point here is that there is natural heterogeneity among individuals in many extraneous variables. Although these extraneous variables are not of any substantive interest, they can potentially have an impact on the response variable. The beauty of a longitudinal study design is that any extraneous factors (regardless of whether they have been measured) that influence the response, and whose influence persists but remains relatively stable throughout the duration of the study (e.g., gender, socioeconomic status, and many genetic, environmental, social, and behavioral factors), are eliminated or blocked out when an individual's responses are compared at two or more occasions. By eliminating these major sources of variability or "noise" from the estimation of within-individual change, a very precise estimate of change can often be obtained.

In summary, the fundamental objective of a longitudinal analysis is the assessment of within-individual changes in the response and the explanation of systematic differences among individuals in their changes. Given that certain individuals change more (or less) than others, the goal of a longitudinal analysis is to determine whether these individuals have larger or smaller values on selected covariates. Finally, in some longitudinal studies, it may also be of interest to make predictions about how specific individuals change over time. In the latter case, longitudinal studies permit more reliable prediction by borrowing information from all individuals to better predict within-individual change over time for a specific individual.

## **2.3 DEFINING FEATURES OF LONGITUDINAL DATA**

At this point we need to introduce some terminology that will be used throughout the remainder of the book. We also introduce some notation for longitudinal data and highlight the main aspects of longitudinal data that complicate their analysis, namely the correlation among repeated observations obtained on the same individual.

### 2.3.1 Terminology

In a longitudinal study the participants, or, more generally, the units being studied, are referred to as *individuals* or *subjects*. In many, but certainly not all, longitudinal studies, the individuals are human subjects. In other longitudinal studies, the individuals may be animals (e.g., laboratory mice or rats). Depending on the specific context, we use the terms *individuals* and *subjects* interchangeably to refer to the participants in a longitudinal study. As mentioned earlier, in a longitudinal study individuals are measured repeatedly at different *occasions* or *times*. Later we will introduce some notation that can distinguish the responses from different individuals in a longitudinal study as well as the repeated measurements on any particular individual. Thus, adopting the terminology introduced so far, the defining feature of a longitudinal study design is that measurements of the response variable are taken on the same *individuals* at several *occasions*.

The number of repeated observations, and their timing, can vary widely from one longitudinal study to another. For example, a clinical trial designed to examine the efficacy of a new analgesic agent may take repeated measures of a self-reported pain scale at baseline and at the end of six 15-minute intervals. This would result in seven repeated measures that are equally separated in time. On the other hand, an observational study of human growth may take measurements of height and weight at 3-month intervals from birth to age 2 years, followed by yearly observations from infancy through young adulthood. By design, the latter study would result in a sequence of repeated measures of height and weight that are unequally separated in time. In both of these examples, the number and the timing of the repeated measurements are the same for all individuals, regardless of whether the occasions of measurement are equally or unequally distributed throughout the duration of the study. Loosely borrowing statistical terminology from the field of experimental design, we refer to the latter studies as being “balanced” over time; that is, all individuals have the same number of repeated measurements obtained at a common set of occasions.

It is an almost inescapable feature of longitudinal studies in the health sciences, especially those where the repeated measurements extend over a relatively long duration, that some individuals will miss their scheduled visit or date of observation. In some studies this may necessitate that observations be made some time before or after the scheduled time. Consequently the sequence of observation times is no longer common to all individuals in the study due to mistimed measurements. In that case we refer to the data as being “unbalanced” over time; that is, the repeated measurements are not obtained at a common set of occasions. Unbalanced longitudinal designs are commonplace when the longitudinal study involves retrospectively collected data (e.g., longitudinal data obtained from medical record databases). Alternatively, highly unbalanced longitudinal data can arise when it is of interest to define the timings of the measurements relative to some benchmark event that occurs during the follow-up period. For example, in a study examining changes in body fat in girls before and after menarche (to be discussed in Section 8.8), the study was designed to begin annual follow-up measurements of body fat prior to menarche and continue for four years after menarche. Although this study design is balanced if the timing of measurements is defined as the time since the baseline measurement, the data are inherently unbalanced if the timing of measurements is defined as the time since an individual experienced menarche. Thus longitudinal studies that are balanced over time when the timing of measurements is defined according to one origin can become highly unbalanced when time is defined in terms of a different origin.

Although longitudinal designs that are unbalanced over time often arise due to happenstance, they are sometimes planned by the investigators. In a “rotating panel” study design, which is commonly used in health surveys to reduce response burden, individuals rotate in and out of the study after providing a pre-determined number of repeated measures. For example, two or more “panels” of individuals are measured repeatedly for a restricted number of occasions, with the first measurement for each “panel” of individuals being staggered. Thus some individuals rotate out (either temporarily or permanently) of the sample, whereas other individuals rotate into the sample. The primary motivation for this type of study design is to reduce costs and the overall burden of participating in the study for any individual, while providing observations at every occasion for some pre-

determined proportion of the sample. An important characteristic of the rotating panel design is that the number and timing of the measurements is pre-determined and by design. Furthermore the decision about whether to obtain a measurement on an individual at any specific occasion is pre-determined a priori by the investigators and is not related to the response variable.

Missing data are a common and challenging problem in longitudinal studies. Indeed, missing data are the rule, not the exception, in longitudinal studies in the health sciences. For example, study participants do not always appear for a scheduled observation, or they may simply leave the study before its completion. When some observations are missing, the data are necessarily unbalanced over time, since not all individuals have the same number of repeated measurements obtained at a common set of occasions. However, to distinguish missing data in a longitudinal study from other kinds of unbalanced data, such data sets are often referred to as being “incomplete.” This distinction is important and emphasizes the fact that an intended measurement on an individual could not be obtained.

One of the consequences of lack of balance and/or missing data is that it requires some care to recover within-individual change. For example, consider a setting where each individual is measured on each of  $n$  occasions. Then consider plotting the mean response at each occasion. Differences in the mean response over time measure the within-individual change. This is because the difference in the means is also the mean of the differences when each subject is measured at every occasion. When data are missing, and especially when there is attrition of subjects whose responses are different from those who remain in the study, then a plot of the mean response over time can be misleading; changes over time may reflect the pattern of missingness or the attrition, and not within-individual change. As we will discuss in later chapters, one will need to examine assumptions and the appropriateness of the analysis carefully to determine the validity of the inferences with unbalanced designs and/or missing data. Although the methods discussed in this book are designed to handle unbalanced designs and missing data, it is worth keeping in mind that it is always preferable to have balanced designs, because these designs can only capture within-individual change.

When longitudinal data are incomplete, there are ramifications for their analysis that go beyond whether a particular statistical method can handle unbalanced longitudinal data. First, when there are missing data, it should be intuitively clear that there must necessarily be some loss of information. Thus there is a price to be paid in terms of efficiency or the precision with which changes over time can be estimated. However, besides causing inefficiency, in some circumstances missing data can introduce bias in the estimates of change. As a result, when longitudinal data are incomplete, the reasons for any missingness must be carefully considered. In Chapters 17 and 18 we discuss some of the consequences of incomplete data in longitudinal studies. In all subsequent chapters we allow for missing data but implicitly make assumptions about the reasons for any missingness. These assumptions are discussed in Section 4.3 and spelled out in greater detail in Chapter 17.

In summary, longitudinal data can be balanced and complete when all individuals are measured at a common set of occasions and there are no missing data. In our experience, longitudinal data in the health sciences are rarely balanced and complete unless the subjects lack human volition (e.g., laboratory rats) or the length of the study is relatively short (e.g., a longitudinal study of the efficacy of an analgesic where the repeated measurements can be obtained in a single study visit). It is far more common to have longitudinal data that are unbalanced and/or incomplete. As a result, to be of real practical use, methods for the analysis of longitudinal data must be able to handle data that are unbalanced over time and possibly incomplete.

Finally, an aspect of longitudinal data that features prominently in their statistical analysis is that repeated measures on the same individual are usually positively correlated. As mentioned earlier, correlated observations are a positive feature of longitudinal data because they provide more precise estimates of the rate of change or the effect of covariates on that rate of change than would be obtained from an equal number of independent observations of different individuals. Nevertheless, the correlation among repeated measures violates the fundamental assumption of independence that is the cornerstone of so many standard regression techniques. In later sections we consider the different

sources and nature of the correlation among longitudinal data, and the potential consequences of not accounting for it in the analysis.

## 2.3.2 Notation

Next we introduce some notation that will be used extensively throughout the book. Let  $Y_{ij}$  denote the response variable for the  $i^{th}$  individual ( $i = 1, \dots, N$ ) at the  $j^{th}$  occasion ( $j = 1, \dots, n$ ). If the repeated measures are assumed to be equally separated in time, this notation will be sufficient. Later, however, we will need to refine the notation to handle the case where the repeated measures are unequally separated and unbalanced over time.

In the statistical literature, the usual convention is to denote a random variable by an uppercase letter (e.g.,  $Y_{ij}$  is the response variable for the  $i^{th}$  individual at the  $j^{th}$  occasion) and the realized value of a random variable by the corresponding lowercase letter (e.g.,  $y_{ij}$  denotes the realized value of  $Y_{ij}$ ). For the most part, we adopt this convention throughout the book. However, whenever we deviate from this convention, it should be clear from the context whether we are referring to a random variable or to its realized value. In [Table 2.1](#) we represent the  $n$  observations (or realized values of  $Y_{ij}$ ) on the  $N$  individuals in a two-dimensional array, with rows corresponding to individuals and columns corresponding to the responses at each occasion. Given that we have  $n$  repeated measures of the response variable on the same individual, we can group these into a  $n \times 1$  response vector, denoted by

**Table 2.1** Tabular representation of longitudinal data, with  $n$  repeated observations on  $N$  individuals.

Individual	Occasion				
	1	2	3	...	$n$
1	$y_{11}$	$y_{12}$	$y_{13}$	...	$y_{1n}$
2	$y_{21}$	$y_{22}$	$y_{23}$	...	$y_{2n}$
.	.	.	.	...	.
.	.	.	.	...	.
.	.	.	.	...	.
$N$	$y_{N1}$	$y_{N2}$	$y_{N3}$	...	$y_{Nn}$

$$Y_i = \begin{pmatrix} Y_{i1} \\ Y_{i2} \\ \vdots \\ Y_{in} \end{pmatrix}.$$

For notational convenience, we can denote the response vectors  $Y_i$  in a completely equivalent way as

$$Y_i = (Y_{i1}, Y_{i2}, \dots, Y_{in})'.$$

(Readers unfamiliar with vectors and matrices should take this opportunity to review the *Gentle Introduction to Vectors and Matrices* in Appendix A; because vectors and matrices are used extensively throughout this book, to simplify notation, the reader is required to have some basic facility with the addition and multiplication of vectors and matrices.<sup>1</sup>)

In the analysis of data from a longitudinal study, the main interest is in the mean response, in particular, changes in the mean response over time and how these changes depend on covariates (e.g., treatment group, exposures). We denote the mean or expectation of each response  $Y_{ij}$  by

$$\mu_j = E(Y_{ij}),$$

where  $E(\cdot)$  can be loosely thought of as denoting a long-run average over a large population of subjects at the  $j^{th}$  occasion. A somewhat more precise definition of the expectation of  $Y_{ij}$  (and of expectation more generally) is that it is a *weighted* average of all the possible values of  $Y_{ij}$ , with weights being the probabilities of occurrence of each possible value. So far our discussion of the mean of  $Y_{ij}$  has assumed that the mean response can change over time; this is reflected in our use of a single-letter subscript for the mean,  $\mu_j$ . In many longitudinal studies the main goal is to relate changes in the mean response over time to covariates. To additionally allow the mean response and, in particular, changes in the mean response, to vary from individual to individual as a function of individual-level covariates, we require the use of double-letter subscripts,

$$\mu_{ij} = E(Y_{ij}).$$

Here, expectation denotes a long-run average over a large subpopulation of subjects who share similar values of the covariates (e.g., subjects assigned to the active treatment group, unexposed subjects) at the  $j^{\text{th}}$  occasion. We refer to  $\mu_{ij}$  as the *conditional* mean response at the  $j^{\text{th}}$  occasion, where the term *conditional* is used to denote the dependence of the mean on covariates. In this notation, the mean response can change over time (denoted by the dependence of  $\mu_{ij}$  on the subscript  $j$ ) and changes in the mean response can be related to individual-level covariates (denoted by the dependence of  $\mu_{ij}$  on the subscript  $i$ ); in Chapter 3 we introduce additional notation that makes the dependence of the mean on the covariates more transparent. A simple illustration of a model for the mean response that depends on time, and that allows changes in the mean response to also depend on covariates, is presented in Section 2.4. In Chapters 5 and 6 we present a detailed discussion of two broad approaches for modeling changes in the mean response over time and for relating these changes to covariates.

Next we consider the correlation or dependence among the  $n$  responses on the same individual. The notions of dependence and independence have precise meanings in statistics. Specifically, two variables are said to be *independent* if the conditional distribution of one of them does not depend on the other. For example, LDL cholesterol level would be considered independent of gender if the distribution of LDL cholesterol level were the same for males and females. Many standard statistical techniques (e.g., linear regression and analysis of variance for a single, univariate response) make the assumption that the study observations are realizations of random variables that are independent of one another. This assumption will be quite reasonable when the study design calls for one observation to be obtained from each individual and individuals are randomly selected from a larger population. The independence assumption is also justified when the study calls for one observation to be obtained from each individual and individuals are randomly assigned to different treatment conditions. Moreover the assumption of independent observations can often be justified on purely physical or scientific grounds when the responses from distinct individuals in the study are considered to be completely unrelated to each other. That is, the response of one individual neither influences or is influenced by the response of another. However, in the case where more than a single observation is obtained on the *same* individual, the assumption of independent observations is simply untenable. That is, the response of an individual on one occasion is very likely to be predictive of the response of the same individual at a future occasion. For example, an individual with a high LDL cholesterol level on one occasion is very likely to also have a high LDL cholesterol level on the next occasion. Put simply, with repeated observations on the same individual, past responses are predictive of future responses. Moreover, with a quantitative response variable, this dependence among the repeated measures on the same individual can be characterized by their correlation. As mentioned earlier, the correlation among repeated measures is a positive feature of longitudinal data because correlated observations provide more precise estimates of the rate of change or the effect of covariates on that rate of change than would be obtained from an equal number of independent observations of different individuals. As we will see in later chapters, models for longitudinal data put the correlation among repeated measures to good advantage when estimating changes in the response over time.

### 2.3.3 Dependence and Correlation

In Section 2.5 we discuss the different sources of correlation among longitudinal data. Before doing so, we must define the term “correlation.” To simplify the discussion of correlation, we consider a simple longitudinal design that is balanced and complete, with  $n$  repeated measurements of the response variable made at a common set of occasions on  $N$  individuals.

Before we can give a formal definition of correlation we need to introduce the notions of *variance* and *covariance*. If we denote the conditional expectation or mean of  $Y_{ij}$  by

$$\mu_{ij} = E(Y_{ij}),$$

then the conditional variance of  $Y_{ij}$  is defined as

$$\sigma_j^2 = E\{Y_{ij} - E(Y_{ij})\}^2 = E(Y_{ij} - \mu_{ij})^2.$$

While  $\mu_{ij}$  provides a measure of the location of the center of the distribution of  $Y_{ij}$ , the conditional variance provides a measure of the spread or dispersion of the values of  $Y_{ij}$  around their conditional mean. The positive square-root of the conditional variance,  $\sigma_j$ , is known as the conditional *standard deviation*. (Readers unfamiliar with expectations and variances are encouraged to take this opportunity to review the *Properties of Expectations and Variances* in Appendix B.) Note that in our discussion of the variance we have implicitly assumed that it can vary from occasion to occasion (reflected in our use of a single-letter subscript,  $\sigma_j^2$ ). In principle, the variance can also be allowed to depend on individual-level covariates; this would require the use of double-letter subscripts. For ease of exposition we have chosen not to do so; in later chapters we will discuss how the variances can vary not only from one occasion to another, but also as a function of selected covariates.

Next we consider the dependence among the responses in a longitudinal study. The conditional *covariance* between the responses at two different occasions, say  $Y_{ij}$  and  $Y_{ik}$ , is denoted by

$$\sigma_{jk} = E\{(Y_{ij} - \mu_{ij})(Y_{ik} - \mu_{ik})\},$$

and provides a measure of the *linear* dependence between  $Y_{ij}$  and  $Y_{ik}$ , given the covariates. The covariance between  $Y_{ij}$  and  $Y_{ik}$  can take on both positive and negative values. When the covariance is zero, there is no linear dependence between the responses at the two occasions (given the covariates). The magnitude of the covariance depends not only on the degree of dependence between the two variables but also on their units of measurement. Any changes in the measurement scales will result in a change in the value of the covariance. For example, the covariance between body weight and LDL cholesterol level will be different if body weight is measured in kilograms rather than pounds. Of note, the covariance of a variable with itself (e.g., the covariance between  $Y_{ij}$  and  $Y_{ij}$ ) is simply the variance of the variable.

While the sign (positive or negative) of the covariance indicates whether there is positive or negative dependence between the two variables, the magnitude of the covariance is somewhat difficult to interpret without comparison to the underlying variability of the two variables. For example, if  $\sigma_{jk} = 10$ , this information alone indicates that there is dependence between  $Y_{ij}$  and  $Y_{ik}$  (since  $\sigma_{jk} \neq 0$ ) and that the dependence is positive (i.e.,  $Y_{ij}$  increases as  $Y_{ik}$  increases, and vice versa). However, depending on the magnitude of the variances of  $Y_{ij}$  and  $Y_{ik}$ ,  $\sigma_{jk} = 10$  may indicate weak or strong dependence. As a result the covariance alone is not too informative; it must be interpreted relative to the magnitude of the variances of the two variables.

To provide a measure of linear dependence between  $Y_{ij}$  and  $Y_{ik}$  that is in some sense free of the units of measurement (or variability) of the two variables, the correlation is widely used.

The conditional correlation between  $Y_{ij}$  and  $Y_{ik}$  is denoted by

$$\rho_{jk} = \frac{E\{(Y_{ij} - \mu_{ij})(Y_{ik} - \mu_{ik})\}}{\sigma_j \sigma_k},$$

where  $\sigma_j$  and  $\sigma_k$  are the conditional standard deviations of  $Y_{ij}$  and  $Y_{ik}$ , respectively. The correlation, unlike the covariance, is a measure of dependence that is unitless or free of the scales of measurement of  $Y_{ij}$  and  $Y_{ik}$ . This is achieved by dividing each variable by its respective standard

deviation. As a result the correlation between body weight and LDL cholesterol level is the same regardless of whether body weight is measured in kilograms or pounds. This makes it a more readily interpretable measure of linear dependence between two variables. Note that when the covariance is zero, so too is the correlation.

By definition, correlation must take values between  $-1$  and  $1$ . Recall that a correlation of  $1$  or  $-1$  is obtained when there is a perfect linear relationship between the two variables. That is, if pairs of values of  $Y_{ij}$  and  $Y_{ik}$  were plotted as points on a two-dimensional scatterplot (assuming the absence of covariates), the resulting points would lie perfectly along a straight line when  $\rho_{jk} = \pm 1$ . As the points depart from a perfect straight-line relationship, the correlation moves closer to zero. A positive correlation implies that one variable increases as the other variable increases. Although two variables that are statistically independent of one another will necessarily be uncorrelated, variables can be uncorrelated without being independent (since correlation only measures *linear* dependence). Statistical independence is a stronger condition than zero correlation; it implies no dependence whatsoever, that is, no *linear* or *non-linear* dependence between the variables. On the other hand, correlation quantifies the degree to which two variables are related or dependent, provided that the dependence is *linear*.

With longitudinal data the repeated measures on the same individual are anticipated to be positively correlated. When the  $n$  repeated measures are collected into a vector  $Y_i = (Y_{i1}, Y_{i2}, \dots, Y_{in})'$ , we can define the variance-covariance matrix to be the following two-dimensional array of conditional variances and covariances:

$$\begin{aligned}\text{Cov} \begin{pmatrix} Y_{i1} \\ Y_{i2} \\ \vdots \\ Y_{in} \end{pmatrix} &= \begin{pmatrix} \text{Var}(Y_{i1}) & \text{Cov}(Y_{i1}, Y_{i2}) & \cdots & \text{Cov}(Y_{i1}, Y_{in}) \\ \text{Cov}(Y_{i2}, Y_{i1}) & \text{Var}(Y_{i2}) & \cdots & \text{Cov}(Y_{i2}, Y_{in}) \\ \vdots & \vdots & \ddots & \vdots \\ \text{Cov}(Y_{in}, Y_{i1}) & \text{Cov}(Y_{in}, Y_{i2}) & \cdots & \text{Var}(Y_{in}) \end{pmatrix} \\ &= \begin{pmatrix} \sigma_{11} & \sigma_{12} & \cdots & \sigma_{1n} \\ \sigma_{21} & \sigma_{22} & \cdots & \sigma_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ \sigma_{n1} & \sigma_{n2} & \cdots & \sigma_{nn} \end{pmatrix},\end{aligned}$$

where  $\text{Cov}(Y_{ij}, Y_{ik}) = \sigma_{jk}$  (and we have implicitly assumed that the variances and covariances are constant across individuals). Note that there is a symmetry to this matrix in the sense that  $\text{Cov}(Y_{ij}, Y_{ik}) = \sigma_{jk} = \sigma_{kj} = \text{Cov}(Y_{ik}, Y_{ij})$ . Also recall that the covariance of a variable with itself is the variance. Thus we can denote

$$\sigma_{kk} = \text{Cov}(Y_{ik}, Y_{ik}) = \text{Var}(Y_{ik}) = \sigma_k^2.$$

For the remainder of the book, and to avoid any potential confusion, we denote the standard deviation and variance of  $Y_{ik}$  by  $\sigma_k$  and  $\sigma_k^2$ , respectively. Also we often refer to the variance-covariance matrix of  $Y_i$  as the covariance (matrix) of  $Y_i$  or simply  $\text{Cov}(Y_i)$ . Thus

$$\text{Cov}(Y_i) = \begin{pmatrix} \sigma_1^2 & \sigma_{12} & \cdots & \sigma_{1n} \\ \sigma_{21} & \sigma_2^2 & \cdots & \sigma_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ \sigma_{n1} & \sigma_{n2} & \cdots & \sigma_n^2 \end{pmatrix}.$$

We can also define the correlation matrix,  $\text{Corr}(Y_i)$ , in terms of a similar two-dimensional array,

$$\text{Corr}(Y_i) = \begin{pmatrix} 1 & \rho_{12} & \cdots & \rho_{1n} \\ \rho_{21} & 1 & \cdots & \rho_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ \rho_{n1} & \rho_{n2} & \cdots & 1 \end{pmatrix}.$$

This matrix is said to be *symmetric* in the sense that  $\text{Corr}(Y_{ij}, Y_{ik}) = \rho_{jk} = \rho_{kj} = \text{Corr}(Y_{ik}, Y_{ij})$ . The diagonal elements of the matrix are all equal to  $1$ , since they denote the correlation of a variable with itself.

With longitudinal data the usual assumptions for standard regression analysis do not hold. In particular, repeated observations on the same individual are not independent and the variance of the

repeated measurements is not usually constant over the duration of the study (e.g., the variance of baseline measurements is often discernibly smaller than post-baseline measurements). With longitudinal data, heterogeneity of variance over time can be accounted for by allowing the elements on the main diagonal of the covariance matrix to differ. The lack of independence among the repeated measurements is accounted for by allowing the off-diagonal elements of the covariance and correlation matrices to be non-zero. Moreover, with longitudinal data, the correlations are expected to be positive and the sequential nature of longitudinal data implies that there may be a pattern to the correlations. For example, a pair of repeated measures that have been obtained close together in time are expected to be more highly correlated than a pair of repeated measures further separated in time. In general, with longitudinal data the correlation among the repeated measures is expected to decline with increasing time separation. In later chapters of this book we will discuss models for the covariance matrix that attempt to capture this structure or pattern in the correlations and that allow the variances to change over time.

In the following section we consider a simple example to highlight the main objectives of a longitudinal analysis and to reinforce the concepts of covariance and correlation that were introduced earlier.

## 2.4 EXAMPLE: TREATMENT OF LEAD-EXPOSED CHILDREN TRIAL

In this simple illustration we consider data from the *Treatment of Lead-Exposed Children Trial*. The TLC trial was a placebo-controlled, randomized study of succimer in children with blood lead levels of 20 to 44  $\mu\text{g}/\text{dL}$ . Recall that the data consist of four repeated measurements of blood lead levels obtained at baseline (or week 0), week 1, week 4, and week 6 on 100 children who were randomly assigned to chelation treatment with succimer or to placebo. Although there were some minor departures from the measurement schedule (e.g., due to mistimed measurements), for the purposes of illustration we regard these data as arising from a balanced design.

# Objectives of Analysis

In general, the main objective of a longitudinal analysis is to describe changes in the mean response over time, and how these changes are related to covariates of interest. In the TLC trial the investigators were interested in determining whether chelation treatment with succimer reduces blood lead levels over time relative to any changes observed in the placebo group. Although the scientific objective of this study is clear, there are many possible ways to express this question in terms of within-individual changes in blood lead levels. For instance, the null hypothesis of no treatment effect on changes in blood lead levels over time could be expressed as

$$H_0: \mu_j(S) = \mu_j(P), \text{ for all } j = 1, \dots, 4,$$

where the notation  $\mu_j(S)$  and  $\mu_j(P)$  is used to denote the mean response at the  $j^{th}$  occasion in the succimer and placebo groups, respectively. This null hypothesis states that the mean responses *at every time point* coincide or are equal in the two treatment groups. As we mentioned earlier, the regression approach to modeling longitudinal data can be formulated in such a way that certain regression parameters correspond to the scientific question of interest. Here, a regression model for the blood lead level data might include main effects for treatment group and time, in addition to their interaction. The null hypothesis given above can then be expressed in terms of the regression parameters for both the main effect of treatment group and the time by treatment group interaction.

Alternatively, the null hypothesis of no treatment effect on changes in blood lead levels over time could be expressed as

$$H_0: \mu_j(S) - \mu_1(S) = \mu_j(P) - \mu_1(P), \text{ for all } j = 2, \dots, 4.$$

This null hypothesis states that all changes in the mean response from baseline are equal in the two treatments groups. Of note, this second version of the null hypothesis is implied by the first. The second version is somewhat less restrictive in that the treatment groups could have differences in means at baseline but identical changes from baseline over time. As we will see later in this book, there are a variety of ways to handle baseline measurements in the analysis of longitudinal data. Once again, a regression model can be formulated corresponding to this second version of the null hypothesis. Specifically, the null hypothesis can be expressed in terms of the regression parameters for the treatment group by time interaction (in a model that includes main effects for treatment group and time).

Finally, a third possibility is to express the null hypothesis in terms of the rate of decline of blood lead levels in the two treatment groups, where the rate of decline or trajectory over time is defined parametrically (e.g., in terms of the slope of a linear response trajectory). However, before we can express and test this null hypothesis, we need to specify more precisely what we mean by rate of decline. In Chapter 6 we will describe how simple parametric (e.g., linear or quadratic) or semiparametric curves (e.g., piecewise linear) can be used to describe trajectories of the mean response changes over time. From a statistical perspective, expressing the null hypothesis in terms of simple parametric curves can result in tests of treatment effects that have greater statistical power.

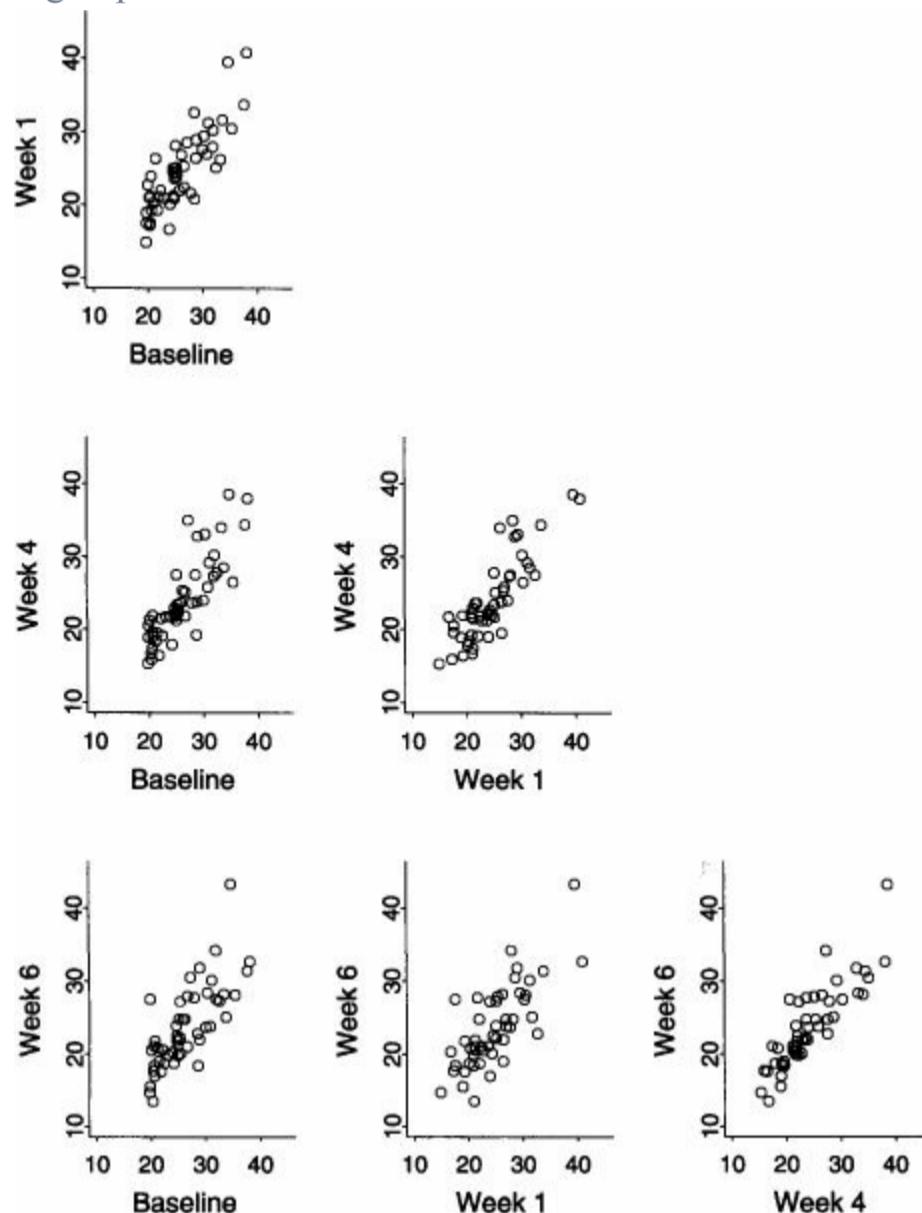
Although there are a number of different ways to express the null hypothesis, in all three instances the main scientific goal is to establish whether changes in the mean response are affected by treatment. More generally, in most longitudinal studies the primary focus is on determining whether the mean response changes over time and whether the changes are related to covariates. The statement of the study hypothesis will depend to a certain extent on the design of the study and the specific goals of the analysis. In Chapters 5 and 6 we will consider this issue in much greater detail.

# Correlation and Covariance

In our discussion of correlation, for ease of exposition, we restrict attention to the longitudinal data from the placebo treated group in this trial. Recall that the data consist of four repeated measurements of blood lead levels obtained at baseline (or week 0), week 1, week 4, and week 6. Thus, for the subset of 50 children who were randomly assigned to the placebo group, we let  $Y_{ij}$  denote the blood lead level for the  $i^{th}$  individual ( $i = 1, \dots, 50$ ) at the  $j^{th}$  occasion ( $j = 1, \dots, 4$ ).

The inter-dependence (or time-dependence) among the four repeated measures of blood lead level can be examined by constructing a scatterplot of each pair of repeated measures. [Figure 2.1](#) displays scatterplots constructed for all six possible pairings of the four repeated measures. [Figure 2.1](#) indicates that there is a relatively strong positive relationship between repeated measures of blood lead levels over time.

**Fig. 2.1** Pairwise scatterplots of blood lead levels at baseline, week 1, week 4, and week 6 for children in the placebo group of the TLC trial.



The estimated covariances and correlations among the four repeated measures are displayed in [Tables 2.2](#) and [2.3](#). Examination of the main diagonal of the covariance matrix reveals that the variances increase over time. In our experience, increasing variance over time is a very common characteristic of longitudinal data. Thus the changing variance of longitudinal data is another type of nuisance problem that is non-standard in most regression settings. Examination of the correlations in [Table 2.3](#) confirms that the correlations are all positive and that the correlation shows a tendency to decrease with increasing time separation.

**Table 2.2** Estimated covariance matrix for the blood lead levels at baseline, week 1, week 4, and week 6 for children in the placebo group of the TLC trial.

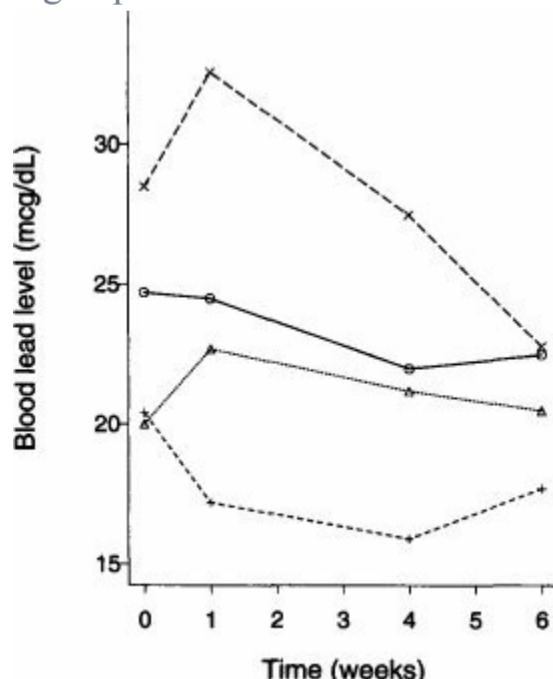
Covariance Matrix			
25.2	22.8	24.3	21.4
22.8	29.8	27.0	23.4
24.3	27.0	33.1	28.2
21.4	23.4	28.2	31.8

**Table 2.3** Estimated correlation matrix for the blood lead levels at baseline, week 1, week 4, and week 6 for children in the placebo group of the TLC trial.

Correlation Matrix			
1.00	0.83	0.84	0.76
0.83	1.00	0.86	0.76
0.84	0.86	1.00	0.87
0.76	0.76	0.87	1.00

While the scatterplots in [Figure 2.1](#) provide a clear indication of the positive correlation among the repeated measures, there is another, albeit less obvious, way to graphically assess the dependence among the repeated measures. This can be achieved using a single scatterplot that plots the responses on the vertical axis and the times of measurements on the horizontal axis, with successive repeated measures on the same individual joined with straight lines; we refer to the resulting display as a *time plot*. The dependence among the repeated measures is assessed by comparing the relative amount of between-subject and within-subject variability. It is usually sufficient, and generally more informative, to produce this scatterplot for only a few randomly selected individuals; it can be very difficult to discern the two distinct sources of variability in a scatterplot based on all of the individuals in the study. In [Figure 2.2](#), based on four randomly selected individuals from the placebo group in the TLC trial, we see that there is very substantial within-subject variability in blood lead levels. This can be discerned from the somewhat jagged appearance of the line segments that join the repeated measures on any individual. In addition there is also substantial between-subject variability. This can be discerned from the fact that some of the individuals have consistently high blood lead levels at all four occasions, while others have consistently low blood lead levels. At first glance this appears to be a very indirect way to assess the degree of dependence among repeated measures and, in our experience, it is not usually the most satisfactory or informative graphical display of that dependence. Nonetheless, it does provide a direct explanation for one of the major sources of the correlation among repeated measures, namely between-individual heterogeneity. In the next section we examine the three major sources of the correlation among repeated measures in a longitudinal study.

**Fig. 2.2** Time plot of blood lead levels at baseline, week 1, week 4, and week 6 for four randomly selected children from the placebo group of the TLC trial.



We have seen that with longitudinal data the usual assumptions for standard regression analysis do not hold. In particular, repeated observations on the same individual are not independent and the variance of the repeated measurements is not usually constant over the duration of the study (e.g., the variance of baseline measurements is often discernibly smaller than post-baseline measurements). The correlation among repeated measures is a positive feature of longitudinal data because correlated observations provide more precise estimates of the rate of change than would be obtained from an equal number of independent observations of different individuals. Although it is important to take this correlation into account in the analysis, the correlations may not be of substantive interest in their own right. If so, we need to accommodate the correlation in an analysis of longitudinal data, but the correlation is not the main focus of the analysis *per se*. Instead, the main interest in any longitudinal study is in describing changes in the mean response over time, and how these changes are related to covariates of interest. For example, in the TLC trial, the main interest is in determining whether chelation treatment with succimer reduces blood lead levels over time relative to any changes in the placebo group. There is no substantive interest in the correlation among the four repeated measurements of blood lead levels obtained at baseline (or week 0), week 1, week 4, and week 6.

## **2.5 SOURCES OF CORRELATION IN LONGITUDINAL DATA**

In this section we consider some of the potential sources of the correlation within longitudinal data. While it is almost an article of faith that longitudinal data are correlated, it is worth pausing to consider why this is the case and, moreover, why longitudinal data are usually positively correlated. Our practical experience with many longitudinal studies in the biological and health sciences has led to the following empirical observations about the nature of the correlation among repeated measures in longitudinal studies: (1) the correlations are positive, (2) the correlations often decrease with increasing time separation, (3) the correlations between repeated measures rarely ever approach zero, even in cases where they are taken many years apart, and (4) the correlation between a pair of repeated measures taken very closely together in time rarely approaches one. These empirical observations have led us to conclude that there are generally three potential sources of variability that have an impact on the correlation among repeated measures on the same individual: (1) between-individual heterogeneity, (2) within-individual biological variation, and (3) measurement error. Next, we examine each of these sources of variability in turn and discuss their impact on the correlation among repeated measures.

# Between-Individual Heterogeneity

The first source of variability is between-subject heterogeneity and this reflects natural variation in individuals' propensity to respond. In any longitudinal study some individuals consistently respond higher than the average, while others consistently respond below the average. Thus one source of the positive correlation among repeated measures is the heterogeneity or variability in the response variable between individuals in the population. For almost every health outcome that might be of interest, we can expect to find some degree of heterogeneity. In a certain sense there are always some individuals who are "high respondents" (e.g., individuals with high blood pressure), some who are "low respondents" (e.g., individuals with low blood pressure), and the remainder who are "medium respondents" (e.g., individuals with blood pressure within the so-called normal range).

The central idea that has been introduced here is that each individual's underlying propensity to respond—whether it be "high," "medium," or "low," and whether it be due to genetic, environmental, social, or behavioral factors (or some combination of these factors)—is shared by all of the repeated measures obtained on that individual. As a result an individual with a high value for the response variable at one occasion will be expected to have a relatively high value at subsequent occasions. Consequently a pair of repeated measures on the same individual will be expected to be more similar than single observations obtained from two randomly selected individuals. That is, part of our intuition for why there is a positive correlation among longitudinal responses is that we expect the repeated responses from the same individual to be more similar than the responses across different individuals.

There can also be heterogeneity among individuals in their response trajectories over time. That is, given a treatment or intervention that should lead to an improvement or increase in the response variable, different individuals will invariably show different gains over time. Changes in the response over time, due to the effects of treatments, interventions, or exposures of some kind, are not expected to be completely uniform across all individuals. There will be some individuals whose gains will be above average, while there will be others whose gains are below average. In cases where there is variability in individuals' response trajectories over time, this can account for not only the positive correlation among repeated measures but also the pattern of decreasing correlations with increasing time separation.

In statistical models for longitudinal data, between-individual variability can be accounted for by the introduction of individual-specific "random effects" (e.g., randomly varying intercepts and slopes). That is, to account for between-individual heterogeneity in propensity to respond, some of the effects or regression coefficients in statistical models are assumed to vary randomly. This topic will be discussed in much greater detail in Chapter 8.

In summary, one important source of variability in longitudinal data that has a direct impact on the correlation among the repeated measures is between-subject variation in the response. Another important source of variability is within-subject variation. The notion here is that even in the absence of any treatment, exposure, or intervention, many health-related outcomes are in a state of so-called dynamic constancy. That is, although an individual's underlying propensity to respond may be "high," and this propensity to respond remains relatively fixed over extended periods of time, the observed sequence of repeated measures on this individual will vary in a random manner around this underlying response level. These random fluctuations can be accounted for by at least two main factors: inherent within-individual biological variation in the response over time and measurement error. Next we examine each of these sources of variability in turn.

# Within-Individual Biological Variation

The inherent biological variability of many health outcomes is an important source of variability that has an impact on the correlation among longitudinal responses. Many health-related variables, for example, blood pressure and self-reported pain, fluctuate considerably even over relatively short intervals of time. These fluctuations may be due to circadian rhythms or perhaps influenced by temperature, light, season, diet, or infection. Of the many health-related variables that change over time, a small number vary in quite predictable cyclical rhythms that may be daily (e.g., body temperature), monthly (e.g., estrogen levels in pre-menopausal women), or seasonal in nature. However, most health-related variables do not have such predictable cyclical rhythms. Instead, a sequence of repeated measures on any particular individual will vary around some homeostatic set point in a random manner. Many of these variables can be thought of as realizations of some biological process or combination of biological processes operating within the individual that vary over time. This variability is sometimes referred to as the *inherent within-individual biological variability*. Inherent biological variability of this kind is evident in almost all measured biological parameters, for example, serum cholesterol, blood pressure, and heart rate.

The notion here is that there is some underlying biological process (or combination of processes) that changes through time in a relatively smooth and continuous fashion. As a result random deviations or departures from an individual's underlying response trajectory are likely to be more similar (e.g., both positive or both negative) when measurements are obtained very close together in time. That is, successive random deviations cannot be assumed to be independent. One consequence of this type of variation is that measurements taken very closely together will typically be more highly correlated than measurements that are further separated in time. That is, all other things being equal, measurements on the same individual will be more alike the closer in time they are taken, and will be less similar the further apart in time. For example, when blood pressure is measured repeatedly at 30-minute intervals, adjacent measurements will be more highly correlated than when the repeated measurements are taken weeks or months apart. Thus inherent within-individual biological variability in the response variable over time introduces serial correlation among repeated measures and results in the correlation matrix having a distinctive structure, with the correlation decreasing as the time separation between repeated measures increases.

Another conceptualization of the within-individual biological variation is in terms of the failure to precisely specify each individual's response trajectory over time. If each individual has a slightly different response trajectory over time, then any misspecification of these response trajectories will induce correlation among the repeated measures. Recall from the definitions of variances and covariances that they are measures of deviations from some model for the mean response. To the extent that the model does not hold for individuals as, for example, when the true trend is quadratic but linear is fitted, the repeated observations will be correlated due to model misspecification. The interdependence between the models for the mean and covariance is a topic that will be discussed at greater length in Chapter 7.

# Measurement Error

A final source of variability in longitudinal data is random measurement error. For some health outcomes, for example, height and weight, variation due to measurement error can be almost negligible (or can be made negligible with the use of more sophisticated measurement instruments). However, for many other outcomes, the variability due to measurement error can be quite substantial. Although this source of variability can account for some of the within-subject variation in many health outcomes, it should not be confused with the inherent (within-individual) biological variability of these outcomes. That is, where it is possible to take two measurements of the response simultaneously on the same individual, thus ruling out the possibility of any inherent biologic variability, the values would not be expected to agree due to the imprecision of the measurement procedure. For example, suppose that the variable of interest is nutrient intake, as determined by a particular biomarker in the blood. Furthermore suppose that a blood sample is drawn on each individual and the vial of blood is divided into two sub-samples that are each subjected to laboratory measurement of the biomarker of interest. In general, these two replicate measures of the biomarker are not expected to agree due to random measurement error.

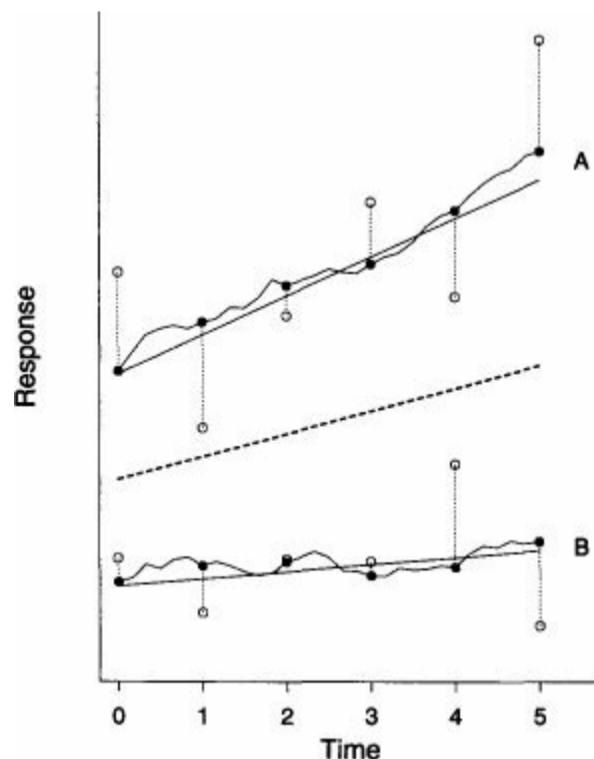
Measurement error is a ubiquitous component of almost all studies, longitudinal or not. Virtually all measurements are fallible, and so it is useful to be able to quantify the relative magnitude of the errors associated with a given measurement procedure. Commonly the precision of the measurement procedure is expressed in terms of a coefficient of *reliability*. The term reliability has a very precise statistical definition and refers to the extent to which replicate measurements, taken under the same conditions, are similar. In the example of the measurement of nutrient intake via a biomarker in the blood, the extent to which these two replicate measures of the biomarker are similar is taken to be an index of the reliability of the biomarker measurement. Note that, at least hypothetically, we could imagine obtaining many such replicate measurements on an individual under as near uniform conditions as possible. In that case an individual's "true" or error-free score is defined as the average of all the (hypothetical) replicate measurements. The statistical definition of reliability then expresses the relative magnitude of the variability of the true scores to the overall variability of the data. That is, reliability is defined as the proportion of the total or overall variability that is due to individual-to-individual variability in the true scores. The reliability of measurements of height is approximately 0.98, while the reliability of measurements of LDL cholesterol can be as low as 0.85. In certain populations, self-reported measures of well-being and quality-of-life can have reliabilities of less than 0.5. In this definition, reliability implicitly depends on the heterogeneity of the true scores in the population of scientific interest. Thus reliability is not a fixed characteristic of the measurement. Because of this, it is preferable to express the precision of a measurement directly in terms of the variance of the measurement errors or alternatively its square-root. (The latter is commonly referred to as the *standard error of measurement*.)

Given that the response variable in most longitudinal studies will be measured with error, what is the potential impact of this variability on the correlation among repeated measures? In general, the effect of unreliability is to "attenuate" or shrink the correlation among the repeated measures closer to zero. For example, if a fallible measure of the response variable has a reliability of 0.8 in the population of interest, the correlation among any pair of repeated measures will be attenuated by a factor of 0.8. In general, the larger the variance of the measurement errors, the greater is the attenuation of the correlation among repeated measures in a longitudinal study. Hence use of a less reliable measurement procedure or instrument will result in repeated measurements with smaller correlations than if a more reliable measurement procedure or instrument had been used.

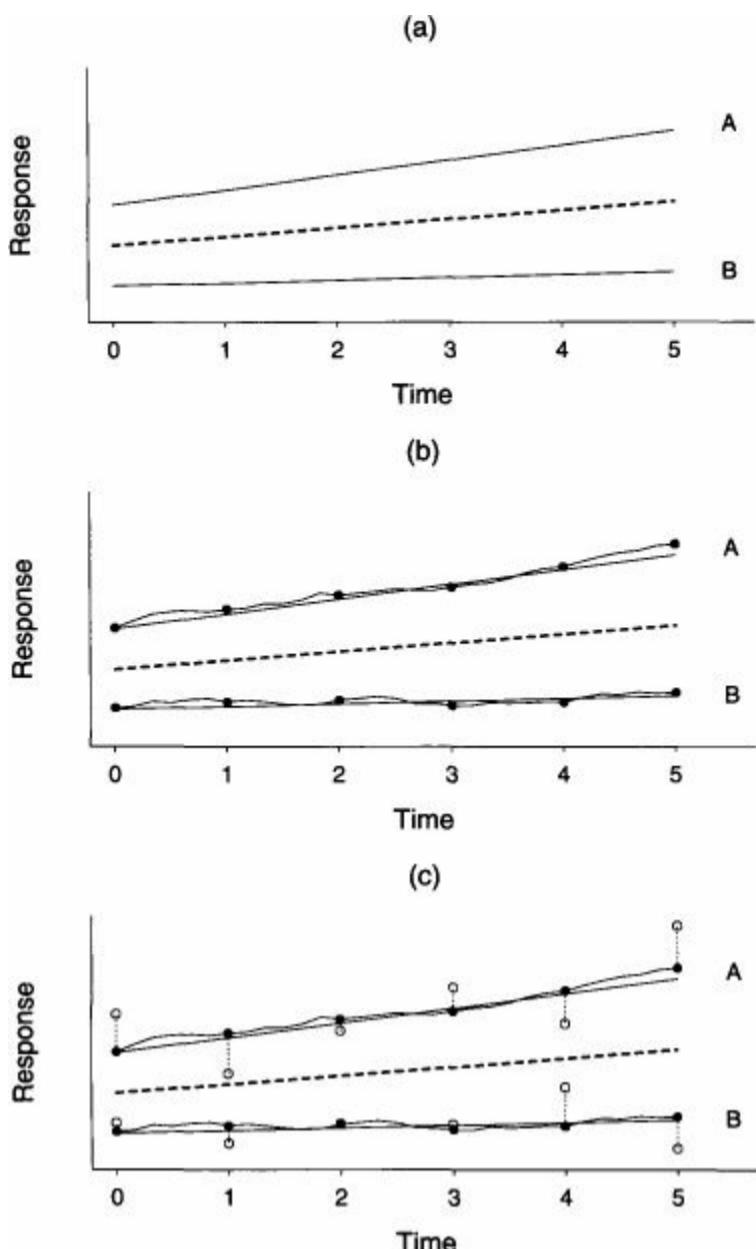
Although we have distinguished two conceptually distinct sources of within-subject variation, within-individual biological variation in the response over time and measurement error, many longitudinal studies will not have sufficient data to estimate these separate sources of variability. That is, for many longitudinal designs it may not be possible to estimate both sources of variability from the data at hand. Instead, for purposes of estimation, both sources may need to be combined into a single component of within-subject variance.

Before concluding our discussion of these three sources of variability in longitudinal data, it is worth pausing to consider the distinctions among them. These three distinct sources of variability can be characterized in a graphical display of longitudinal data on two hypothetical individuals (see [Figure 2.3](#)). [Figure 2.3](#) displays six repeated measurements of the response (denoted by empty circles) on two individuals, say individual A and B. The three sources of variability in the response can be distinguished by considering the additional variability in the response that each source produces. The contributions of these three sources of variability are highlighted in [Figure 2.4](#). In [Figure 2.4\(a\)](#) the between-subject variation is reflected in the degree of separation among the true underlying “response profiles”. These two hypothetical response profiles can be thought of as being representative of the response trajectory for “high” and “low” respondents. If between-subject heterogeneity were the only source of variability in longitudinal data, then the six repeated measures on these two hypothetical individuals would fall along the corresponding response profiles in [Figure 2.4\(a\)](#). However, in addition to between-individual variation, there is within-individual variation. [Figure 2.4\(b\)](#) illustrates how a sequence of repeated measures on any individual might vary in a random manner around their long-run average (or “true” underlying response) due to inherent within-individual biological variation in the response over time. In the absence of any measurement error, repeated measures on these two hypothetical individuals would fall along the corresponding jagged curves; these error-free repeated measures are denoted by the solid circles. However, because of the imprecision of the measurement procedure, the actual repeated measures on these two hypothetical individuals (denoted by empty circles) vary in a random manner around the corresponding jagged lines (see [Figure 2.4\(c\)](#)). The relative magnitude of the between-individual and within-individual sources of variability will be different from one health outcome to another. Their relative magnitude is an important determinant of the degree of correlation among repeated measures.

**Fig. 2.3** Graphical representation of the three sources of variability in longitudinal data for two hypothetical individuals: • denotes repeated measure free of measurement error, ◊ denotes observed repeated measure with measurement error.



**Fig. 2.4** Graphical representation of the cumulative impact of three sources of variability in longitudinal data: (a) between-individual heterogeneity, (b) within-individual biological variation (where • denotes repeated measure free of measurement error), and (c) measurement error (where ◊ denotes observed repeated measure with measurement error).



Finally, we consider the impact of these three sources of variability on the correlation among repeated measures and briefly discuss the potential consequences of ignoring the correlation. Earlier we described four empirical observations about the correlation among repeated measures in longitudinal studies. Here we consider how the three sources of variability in longitudinal data can account for these empirical observations. First, we noted that the correlations among repeated measures are positive. The positive correlation among repeated measures is a direct consequence of both between-individual heterogeneity and within-individual biological variation in the response over time. These two sources of variability act in union to induce positive correlation among the repeated measures. Second, we noted that the correlation tends to decrease with increasing time separation. This is a direct consequence of the inherent within-individual biological variation in the response over time and/or between-individual heterogeneity of response trajectories over time. Third, it was noted that the correlations between repeated measures rarely approach zero, even in cases where the repeated measures are taken many years apart. This is a direct consequence of between-subject heterogeneity in the underlying propensity to respond. That is, an individual's propensity to respond persists across all repeated measures on that individual, regardless of how far apart the measurements are in time. Finally, we noted that the correlation between a pair of repeated measures taken very closely together in time rarely approaches one. This final observation is a direct consequence of measurement error. The correlation between any pair of repeated measurements, regardless of how close the measurement occasions, is constrained by the reliability of the measurement procedure.

While it is likely that all three sources of variability contribute to the variability of longitudinal data, one may be more dominant than another. It may not always be necessary, or indeed possible, to separately estimate these three unique sources of variability. This issue will be examined more closely in Chapter 8. Finally, we remind the reader of the definitions of variance and covariance given in Section 2.3. The conditional variance and covariance are measures defined in terms of a particular model for the conditional mean response over time. As a result there is a subtle

interdependence between the model for the mean response and the model for the covariance. To the extent that the model for the mean response does not fit the data well, the observations will be correlated and overdispersed due to misspecification of the model for the mean response. The interdependence between the models for the mean and covariance, and the ramifications of this interdependence for model selection, are discussed in greater depth in Chapter 7.

# Consequences of Ignoring Correlation among Longitudinal Data

We have seen that longitudinal data are usually positively correlated, and that the strength of the correlation is often a decreasing function of the time separation. Next we consider the potential implications of ignoring the correlation among the repeated measures. In later chapters of this book we will discuss this topic in greater detail. Here we provide a hint of the potential impact of ignoring the correlation with a simple illustration using data from the *Treatment of Lead-Exposed Children Trial*. Consider only the first two repeated measures from this study, taken at baseline (or week 0) and week 1. Suppose that it is of interest to determine whether there has been a change in the mean response over time. A very natural estimate of the change in the mean response over time is

$$\hat{\delta} = \hat{\mu}_2 - \hat{\mu}_1,$$

where

$$\hat{\mu}_j = \frac{1}{N} \sum_{i=1}^N Y_{ij}.$$

For the data from the TLC trial (see [Table 1.2](#)), the estimate of the change in the mean response over time in the succimer group is  $-13.0$  (or  $13.5 - 26.5$ ). Of course, this estimate is not of much use without some estimate of its sampling variability. To obtain the standard error (SE), we need to estimate the variability of this estimator of change. An expression for the variance of  $\hat{\delta}$  is given by

$$\text{Var}(\hat{\delta}) = \text{Var} \left\{ \frac{1}{N} \sum_{i=1}^N (Y_{i2} - Y_{i1}) \right\} = \frac{1}{N} (\sigma_1^2 + \sigma_2^2 - 2\sigma_{12}).$$

It is the inclusion of the last term,  $-2\sigma_{12}$ , in the expression above that accounts for the correlation among the first two repeated measures. For the data at hand, we can substitute estimates of the variances ( $\hat{\sigma}_1^2 = 25.2$  and  $\hat{\sigma}_2^2 = 58.9$ ) and covariance ( $\hat{\sigma}_{12} = 15.5$ ) for the succimer group into this expression to obtain the following estimate of the variance of  $\hat{\delta}$ :

$$\widehat{\text{Var}}(\hat{\delta}) = \frac{1}{50} \{25.2 + 58.9 - 2(15.5)\} = 1.06.$$

If we had simply ignored the fact that the data are correlated and proceeded with an analysis assuming that all observations are independent (and hence uncorrelated, with zero covariance), we would instead have obtained the following (incorrect) estimate of the variance of  $\hat{\delta}$ ,

$$\frac{1}{50} (25.2 + 58.9) = 1.68,$$

which is approximately 1.6 times larger. Thus, in this very simple illustration, ignoring the correlation leads to quite discernible overestimation of the variability of the estimate of change. This in turn would lead to an overly pessimistic estimate of precision, resulting in standard errors that are too large, confidence intervals that are too wide, and  $p$ -values for the test of  $H_0: \delta = 0$  that are too large. In summary, failure to take account of the correlation among the repeated measures will, in general, result in incorrect estimates of the sampling variability, which can lead to quite misleading inferences. This topic will be discussed at much greater length in later chapters.

## 2.6 FURTHER READING

A compelling illustration of the strengths of a longitudinal study design can be found in Chapter 1, Section 1.1, of Diggle *et al.* (2002).

### Problems

**2.1** The *Treatment of Lead-Exposed Children* (TLC) trial was a placebo-controlled, randomized study of succimer (a chelating agent) in children with blood lead levels of 20 to 44 micrograms/dL. Recall that the data consist of four repeated measurements of blood lead levels obtained at baseline (or week 0), week 1, week 4, and week 6 on 100 children who were randomly assigned to chelation treatment with succimer or to placebo. For this problem set we focus only on the 50 children assigned to chelation treatment with succimer.

The raw data are stored in an external file: `lead.dat`

Each row of the data set contains the following 5 variables:

ID Y<sub>1</sub> Y<sub>2</sub> Y<sub>3</sub> Y<sub>4</sub>

**2.1.1** Read the data from the external file and calculate the sample means, standard deviations, and variances of the blood lead levels at each occasion.

**2.1.2** Construct a time plot of the mean blood lead levels versus time (in weeks). Describe the general characteristics of the time trend.

**2.1.3** Calculate the  $4 \times 4$  covariance and correlation matrices for the four repeated measures of blood lead levels.

**2.1.4** Verify that the diagonal elements of the covariance matrix are the variances by comparing to the descriptive statistics obtained in Problem 2.1.1.

**2.1.5** Verify that the correlation between blood lead levels at baseline (week 0) and week 1 is equal to the covariance between blood lead levels at baseline and week 1, divided by the product of the standard deviations of the blood lead levels at baseline and week 1.

<sup>1</sup> Another common convention in the statistical literature is the use of bold type for vectors (and sometimes for matrices). As it will be clear from the context, we do not do so throughout this book.

## *Part II*

# *Linear Models for Longitudinal Continuous Data*

# *Chapter 3*

## *Overview of Linear Models for Longitudinal Data*

### **3.1 INTRODUCTION**

In Part II the focus is exclusively on linear models for longitudinal data with response variables that are continuous and have distributions that are approximately symmetric, without excessively long tails (or skewness) or outliers. The models for longitudinal data presented in Part II also provide the foundations for more general models for longitudinal data when the response variable is discrete or a count. In this chapter we introduce some vector and matrix notation and present a general linear regression model for longitudinal data. A specific feature of the model is that the mean response is linear in the regression parameters. We present a broad overview of different approaches for modeling the mean response over time and for accounting for the correlation among repeated measures on the same individual. These topics are discussed in much greater depth in subsequent chapters. We also consider some elementary descriptive methods for exploring longitudinal data, especially trends in the mean response over time. We conclude the chapter with an historical survey of some of the earliest developments in methods for analyzing longitudinal and repeated measures data.

We must emphasize at the outset that the statistical methods presented in Part II use the assumption that the longitudinal responses have an approximate multivariate normal distribution to derive estimates and statistical tests, but do not require it. That is, the methods discussed in Part II are based on, but do not require, the assumption that the responses have a multivariate normal distribution. Given this distributional assumption, the method of maximum likelihood, presented in Chapter 4, provides a very general technique for estimation and for inference. In Chapter 4 we briefly discuss statistical methods for constructing confidence intervals, testing hypotheses, and assessing the adequacy of models. In Chapter 13, where we discuss alternative methods of estimation, it will become more apparent that we do not require the assumption of multivariate normality.

## **3.2 NOTATION AND DISTRIBUTIONAL ASSUMPTIONS**

In this section we introduce some vector and matrix notation that will be used extensively throughout the remainder of the book. Readers without any prior exposure to matrix algebra are encouraged to review the introduction to vectors and matrices presented in Appendix A; we guarantee that the small investment involved in mastering the material in Appendix A will pay handsome dividends later. Throughout this book we do not presume that the reader has a profound understanding of matrix algebra; however, some basic facility with the addition and multiplication of vectors and matrices is required. As will soon become apparent, our primary motivation for the use of vectors and matrices is the compactness with which multivariate statistical techniques can be presented and described when expressed in vector and matrix notation.

# Notation

In Chapter 2 we assumed that a sample of  $N$  subjects are measured repeatedly over time. We let  $Y_{ij}$  denote the response variable for the  $i^{th}$  subject on the  $j^{th}$  measurement occasion. As was mentioned in Chapter 2, either by design or happenstance, subjects may not have the same number of repeated measures and may not be measured at the same set of occasions. To accommodate both of these features, we assume that there are  $n_i$  repeated measurements of the response on the  $i^{th}$  subject and that each  $Y_{ij}$  is observed at time  $t_{ij}$ . For example, a study may be designed to take repeated measurements on all subjects at the same set of  $n$  occasions. However, missing data are a common problem in almost all longitudinal studies, and some subjects may not have observations at all  $n$  occasions (i.e.,  $n_i$  denotes the number of *observed* responses on the  $i^{th}$  subject, where  $n_i \leq n$ ). Missing data not only produce a varying number of repeated measurements of subjects in a longitudinal study but also have important consequences for the validity of any method of analysis. In Section 4.3 we outline some of the key issues and assumptions required for valid analyses when there are missing data; this topic is discussed in greater detail in Chapter 17. In addition to missing data, there may be mistimed measurements, in the sense that measurements are not obtained at the planned  $n$  occasions; instead, they are obtained some time before or after the intended measurement occasions. Thus both the number and the timing of the repeated measurements may not be common for all subjects. In later chapters the times of measurement,  $t_{ij}$ , are used to model trends in the mean response; they may also be required to appropriately account for the covariance among the repeated measurements.

It is convenient to group the  $n_i$  repeated measures of the response variable for the  $i^{th}$  subject into an  $n_i \times 1$  vector

$$Y_i = \begin{pmatrix} Y_{i1} \\ Y_{i2} \\ \vdots \\ Y_{in_i} \end{pmatrix}, \quad i = 1, \dots, N.$$

Note that the vector  $Y_i$  is simply a time-ordered collection of the  $n_i$  response variables for the  $i^{th}$  subject. The  $Y_{ij}$ 's are often called the components, entries, or elements of  $Y_i$ . The vector  $Y_i$  is said to be of *order*  $n_i \times 1$ , meaning that it consists of  $n_i$  rows and 1 column of elements.

The vectors of responses,  $Y_i$ , for the  $N$  subjects are assumed to be independent of one another. Note, however, that while the vectors of responses obtained on different subjects can usually be assumed to be independent of one another (e.g., repeated measures of a health outcome for one patient in a clinical trial are not expected to predict or influence the health outcomes for another patient in the same trial), the repeated measures on the same subject are emphatically not assumed to be independent observations.

When the number of repeated measures is the same for all subjects in the study (and there are no missing data), it is not necessary to include the index  $i$  in  $n_i$  (since  $n_i = n$  for  $i = 1, \dots, N$ ). Similarly, if the repeated measures are observed at the same set of occasions, it is not necessary to include the index  $i$  in  $t_{ij}$  (since  $t_{ij} = t_j$  for  $i = 1, \dots, N$ ). For example, in the *Treatment of Lead-Exposed Children Trial* all subjects had the same number of repeated measures,  $n = 4$ , and were measured at the same set of occasions,  $\{t_1 = 0, t_2 = 1, t_3 = 4, t_4 = 6\}$ .

Associated with each response,  $Y_{ij}$ , there is a  $p \times 1$  vector of covariates

$$X_{ij} = \begin{pmatrix} X_{ij1} \\ X_{ij2} \\ \vdots \\ X_{ijp} \end{pmatrix}, \quad i = 1, \dots, N; \quad j = 1, \dots, n_i.$$

Note that  $X_{ij}$  is a vector of covariates associated with  $Y_{ij}$ , the response variable for the  $i^{th}$  subject at the  $j^{th}$  occasion. The  $p$  rows of  $X_{ij}$  correspond to different covariates. There is a corresponding vector of covariates associated with each of the  $n_i$  repeated measurements on the  $i^{th}$  subject. That is,

$X_{i1}$  is a  $p \times 1$  vector whose elements are the covariate values associated with the response variable for the  $i^{th}$  subject at the 1<sup>st</sup> measurement occasion,  $X_{i2}$  is a  $p \times 1$  vector whose elements are the covariate values associated with the response variable for the  $i^{th}$  subject at the 2<sup>nd</sup> measurement occasion, and so on. The vector  $X_{ij}$  may include two main types of covariates: covariates whose values do not change throughout the duration of the study and covariates whose values change over time. Examples of the former include gender and fixed experimental treatments. Examples of the latter include time since baseline, current smoking status, and environmental exposures. In the former case, the same values of the covariates are replicated in the corresponding rows of  $X_{ij}$  for  $j = 1, \dots, n_i$ . In the latter case, the values taken by the covariates can vary over time (for at least some individuals) and the values in the corresponding rows of  $X_{ij}$  can be different at each measurement occasion. The inclusion of time-varying covariates whose values at any occasion cannot be predicted (e.g., current smoking status) can raise subtle issues concerning the interpretation and estimation of the resulting models. A discussion of these issues is deferred until Chapter 13 (see Section 13.5).

We can group the vectors of covariates into an  $n_i \times p$  matrix of covariates:

$$X_i = \begin{pmatrix} X'_{i1} \\ X'_{i2} \\ \vdots \\ X'_{in_i} \end{pmatrix}, \quad i = 1, \dots, N,$$

where  $X'_{ij}$  denotes the *transpose* of the vector of covariates,  $X_{ij}$ . Recall that the *transpose* is a function that interchanges the rows and columns of a matrix (see Appendix A); thus  $X'_{ij}$  denotes a  $1 \times p$  row vector of covariates for the  $i^{th}$  subject at the  $j^{th}$  occasion. The matrix  $X_i$  is simply an ordered collection of the values of the  $p$  covariates for the  $i^{th}$  subject at each of the  $n_i$  measurement occasions. That is,

$$X_i = \begin{pmatrix} X_{i11} & X_{i12} & \cdots & X_{i1p} \\ X_{i21} & X_{i22} & \cdots & X_{i2p} \\ \vdots & \vdots & \ddots & \vdots \\ X_{in_i1} & X_{in_i2} & \cdots & X_{in_ip} \end{pmatrix},$$

where the rows of  $X_i$  correspond to the covariates associated with the responses at the  $n_i$  different measurement occasions, and the columns of  $X_i$  correspond to the  $p$  distinct covariates.

By now it should be apparent that the use of vectors and matrices can greatly facilitate exposition by allowing the repeated measurements on the response variable and the covariates to be expressed in a succinct manner. So far we have assumed that each subject in the study has a vector of repeated responses, denoted by  $Y_i$ , and associated with each repeated measure, a vector of  $p$  covariates which can be collectively grouped into a matrix,  $X_i$ . Later we will present a simple numerical example to reinforce the reader's understanding of the vector and matrix notation used so far.

Next we consider a linear regression model for changes in the mean response over time and for relating the changes to the covariates,

$$(3.1) \quad Y_{ij} = \beta_1 X_{i1j} + \beta_2 X_{i2j} + \cdots + \beta_p X_{ipj} + e_{ij}, \quad j = 1, \dots, n_i;$$

where  $\beta_1, \dots, \beta_p$  are unknown regression coefficients relating the mean of  $Y_{ij}$  to its corresponding covariates. This regression model describes how the responses at every occasion are related to the covariates. That is, there are  $n_i$  separate regression equations for the response variable at each of the  $n_i$  occasions

$$\begin{aligned} Y_{i1} &= \beta_1 X_{i11} + \beta_2 X_{i12} + \cdots + \beta_p X_{i1p} + e_{i1} = X'_{i1}\beta + e_{i1}, \\ Y_{i2} &= \beta_1 X_{i21} + \beta_2 X_{i22} + \cdots + \beta_p X_{i2p} + e_{i2} = X'_{i2}\beta + e_{i2}, \\ &\vdots \quad \vdots \quad \vdots \quad \vdots \quad \vdots \end{aligned}$$

$$(3.2) \quad Y_{in_i} = \beta_1 X_{in_i1} + \beta_2 X_{in_i2} + \cdots + \beta_p X_{in_ip} + e_{in_i} = X'_{in_i}\beta + e_{in_i},$$

where the unknown regression parameters are grouped together into a  $p \times 1$  vector,  $\beta = (\beta_1, \beta_2, \dots, \beta_p)'$ . Here the  $e_{ij}$  are random errors, with mean zero, representing deviations of the responses from

their corresponding predicted means

$$E(Y_{ij}|X_{ij}) = \beta_1 X_{ij1} + \beta_2 X_{ij2} + \cdots + \beta_p X_{ijp}.$$

Typically, although not always,  $X_{ij1} = 1$  for all  $i$  and  $j$ , and then  $\beta_1$  is the intercept term in the model. Our use of  $\beta_1$ , rather than  $\beta_0$  or  $\alpha$ , to denote the intercept is somewhat arbitrary but does lead to minor simplification of the notation used throughout the book.

Finally, using vector and matrix notation, the regression model given by (3.1) or (3.2) can be expressed in an even more compact form,

$$(3.3) \quad Y_i = X_i\beta + e_i,$$

where  $e_i = (e_{i1}, e_{i2}, \dots, e_{in_i})'$  is an  $n_i \times 1$  vector of random errors associated with the corresponding elements of the vector of responses on the  $i^{th}$  subject. The regression model given by (3.3) is simply a shorthand representation for

$$\begin{pmatrix} Y_{i1} \\ Y_{i2} \\ \vdots \\ Y_{in_i} \end{pmatrix} = \begin{pmatrix} X_{i11} & X_{i12} & \cdots & X_{i1p} \\ X_{i21} & X_{i22} & \cdots & X_{i2p} \\ \vdots & \vdots & \ddots & \vdots \\ X_{in_i1} & X_{in_i2} & \cdots & X_{in_ip} \end{pmatrix} \begin{pmatrix} \beta_1 \\ \beta_2 \\ \vdots \\ \beta_p \end{pmatrix} + \begin{pmatrix} e_{i1} \\ e_{i2} \\ \vdots \\ e_{in_i} \end{pmatrix}.$$

By comparing (3.2) and (3.3), it should now be apparent to the reader that one of the chief advantages of using vector and matrix notation is that regression models relating longitudinal responses to multiple predictors can be expressed in a very economical fashion.

Thus far we have made no assumption about the conditional distribution of  $Y_i$ , given the covariates. The only assumption made is that the mean of the longitudinal response vector is related to the covariates via the linear regression model given above. Before discussing assumptions about the conditional distribution of  $Y_i$ , let us return to the *Treatment of Lead-Exposed Children Trial* in order to reinforce understanding of the notation introduced so far and to clarify how the regression parameters in (3.3) describe patterns of change in the mean response and their relation to covariates.

# Illustration: Treatment of Lead-Exposed Children Trial

Recall that in the *Treatment of Lead-Exposed Children Trial* there are 100 study participants who have blood lead levels measured at the same set of four occasions: baseline (or week 0), week 1, week 4, and week 6. Since all subjects have the same number of repeated measures observed at the same set of occasions, the index  $i$  can be dropped from both  $n_i$  and  $t_{ij}$ . That is,  $n_1 = n_2 = \dots = n_N = n$ , and similarly  $t_{1j} = t_{2j} = \dots = t_{Nj} = t_j$  for  $j = 1, \dots, 4$ . In the TLC trial the response vector is of length 4 ( $n = 4$ ), and all subjects are measured at the same set of occasions:  $t_1 = 0$ ,  $t_2 = 1$ ,  $t_3 = 4$ , and  $t_4 = 6$ .

Next suppose that it is of interest to fit a model to the mean response that assumes that the mean blood lead level changes linearly over time, but at a rate that might be different for the two treatment groups. In particular, we might want to fit a model where the two treatment groups have the same intercept (or mean response at baseline) but different slopes. This can be represented in the following regression model

$$\begin{aligned} Y_{ij} &= \beta_1 X_{ij1} + \beta_2 X_{ij2} + \beta_3 X_{ij3} + e_{ij} \\ &= X'_{ij}\beta + e_{ij}, \end{aligned}$$

where  $X_{ij1} = 1$  for all  $i$  and all  $j$ . That is,  $X_{ij1} = 1$  for all subjects and at all measurement occasions, and thus  $\beta_1$  is an intercept term. The second covariate,  $X_{ij2} = t_j$ , represents the week in which the blood lead level was obtained. Finally,  $X_{ij3} = t_j \times \text{Group}_i$  where  $\text{Group}_i = 1$  if the  $i^{\text{th}}$  subject is assigned to the succimer group and  $\text{Group}_i = 0$  if the  $i^{\text{th}}$  subject is assigned to the placebo group. As we will show, this coding of  $X_{ij2}$  and  $X_{ij3}$  allows the slopes for time to differ for the two treatment groups. The three covariates can be grouped into a  $3 \times 1$  vector of covariates  $X_{ij}$ . Thus for children in the placebo group

$$E(Y_{ij}|X_{ij}) = \beta_1 + \beta_2 t_j,$$

where  $\beta_1$  represents the mean blood lead level at baseline (week = 0) and  $\beta_2$  has interpretation as the change in mean blood level (in  $\mu\text{g}/\text{dL}$ ) per week. For example, the expected change in mean blood level, from baseline to 6 weeks, is  $\beta_2 \times 6$  for children in the placebo group. Similarly for children in the succimer group

$$E(Y_{ij}|X_{ij}) = \beta_1 + (\beta_2 + \beta_3)t_j,$$

where  $\beta_1$  represents the mean blood level at baseline (assumed to be the same as in the placebo group since the trial randomized subjects to the two groups) and  $\beta_2 + \beta_3$  has interpretation as the change in mean blood level per week. Thus, if the two treatment groups differ in their rates of decline in blood lead levels, then  $\beta_3 \neq 0$ . The regression parameters have useful interpretations that bear directly on questions of scientific interest. Moreover hypotheses of interest can be expressed in terms of the absence (or setting to zero) of certain regression parameters. For example, the hypothesis that the two treatments are equally effective in reducing blood lead levels corresponds to a hypothesis that  $\beta_3 = 0$ .

To reinforce the vector and matrix notation we have introduced in this section, it is instructive to examine the matrix of covariates,  $X_i$ , and the realized values of  $Y_i$ , for any particular subject in the trial. For example, for the study participant with ID = 79 (see [Table 1.1](#)) the realized values of  $Y_i$  are

$$\begin{pmatrix} 30.8 \\ 26.9 \\ 25.8 \\ 23.8 \end{pmatrix}.$$

This individual<sup>1</sup> was assigned to treatment with placebo and thus has the following matrix of covariates,  $X_i$ .

$$\begin{pmatrix} 1 & 0 & 0 \\ 1 & 1 & 0 \\ 1 & 4 & 0 \\ 1 & 6 & 0 \end{pmatrix};$$

the latter is often referred to as the *design matrix*. The four rows of  $X_i$  correspond to the covariates associated with the blood lead levels at the four measurement occasions (weeks 0, 1, 4, and 6). The elements of the first column are all ones (and multiply the intercept term,  $\beta_1$ ). The second column contains values that denote the week in which the blood lead level was obtained. All the elements of the third column are zero (for subjects assigned to the placebo group). On the other hand, for the study participant with ID = 8, the realized values of  $Y_i$  are

$$\begin{pmatrix} 26.5 \\ 14.8 \\ 19.5 \\ 21.0 \end{pmatrix}.$$

This individual was assigned to treatment with succimer and thus has the following design matrix or matrix of covariates,  $X_i$ :

$$\begin{pmatrix} 1 & 0 & 0 \\ 1 & 1 & 1 \\ 1 & 4 & 4 \\ 1 & 6 & 6 \end{pmatrix}.$$

Finally, using vectors and matrices, the model for the mean blood lead levels can be represented as

$$E(Y_i|X_i) = X_i\beta,$$

where

$$E(Y_i|X_i) = \begin{pmatrix} E(Y_{i1}|X_{i1}) \\ E(Y_{i2}|X_{i2}) \\ E(Y_{i3}|X_{i3}) \\ E(Y_{i4}|X_{i4}) \end{pmatrix} = \begin{pmatrix} 1 & 0 & 0 \\ 1 & 1 & 0 \\ 1 & 4 & 0 \\ 1 & 6 & 0 \end{pmatrix} \begin{pmatrix} \beta_1 \\ \beta_2 \\ \beta_3 \end{pmatrix} = \begin{pmatrix} \beta_1 \\ \beta_1 + \beta_2 \\ \beta_1 + 4\beta_2 \\ \beta_1 + 6\beta_2 \end{pmatrix}$$

for children in the placebo group, and

$$E(Y_i|X_i) = \begin{pmatrix} E(Y_{i1}|X_{i1}) \\ E(Y_{i2}|X_{i2}) \\ E(Y_{i3}|X_{i3}) \\ E(Y_{i4}|X_{i4}) \end{pmatrix} = \begin{pmatrix} 1 & 0 & 0 \\ 1 & 1 & 1 \\ 1 & 4 & 4 \\ 1 & 6 & 6 \end{pmatrix} \begin{pmatrix} \beta_1 \\ \beta_2 \\ \beta_3 \end{pmatrix} = \begin{pmatrix} \beta_1 \\ \beta_1 + (\beta_2 + \beta_3) \\ \beta_1 + 4(\beta_2 + \beta_3) \\ \beta_1 + 6(\beta_2 + \beta_3) \end{pmatrix}$$

for children in the succimer group.

# Distributional Assumptions

So far the only assumptions made concern patterns of change in the mean response over time and their relation to covariates. Specifically, given that the vector of random errors,  $e_i$ , is assumed to have mean zero, the regression model given by (3.3) implies that

$$(3.4) E(Y_i|X_i) = \mu_i = X_i\beta,$$

where  $\mu_i = (\mu_{i1}, \dots, \mu_{in_i})'$  is the  $n_i \times 1$  vector of conditional means for the  $i^{th}$  individual, with  $\mu_{ij} = E(Y_{ij}|X_i) = E(Y_{ij}|X_{ij})$ . This model describes how the vector of mean responses is related to the covariates.

Next we consider distributional assumptions concerning the vector of random errors,  $e_i$ . The response vector  $Y_i$  in (3.3) is assumed to be composed of two components, a “systematic component,”  $X_i\beta$ , and a random component,  $e_i$ . The systematic component implies that the mean response can be expressed as a simple weighted sum of the fixed, but unknown, regression coefficients,  $\beta$ . The random variability of  $Y_i$  arises from the addition of  $e_i$ , the “random component.” This implies that assumptions made about the shape of the distribution of the random errors translate into assumptions about the shape of the conditional distribution of  $Y_i$  given  $X_i$ . Thus, in a certain sense, we can almost interchangeably refer to the distribution of either the errors,  $e_i$ , or the responses,  $Y_i$ , their respective distributions differ only in terms of a shift in location. That is, the errors have a distribution with a mean that is centered at zero, while the conditional distribution of  $Y_i$  given  $X_i$  is of the same form except that the mean is centered at  $X_i\beta$ . As a result throughout this book we will interchangeably refer to the distributions of  $Y_i$  and  $e_i$  and, more specifically, the covariance matrix of  $Y_i$  and  $e_i$ . Note that in discussing the distribution of  $Y_i$ , it should be understood that we are always referring to the *conditional* distribution of  $Y_i$  given the covariates,  $X_i$ .

Next  $Y_i$ , the vector of continuous responses, is assumed to have a conditional distribution that is multivariate normal, with mean response vector

$$E(Y_i|X_i) = \mu_i = X_i\beta,$$

and covariance matrix

$$\Sigma_i = \text{Cov}(Y_i|X_i).$$

The multivariate normal distribution is completely specified by the vector of means,  $\mu_i$ , and the covariance matrix,  $\Sigma_i$ . The multivariate normal distribution can be considered to be the multivariate analogue of the univariate normal distribution. Indeed, if  $Y_i$  has a conditional distribution that is multivariate normal, then each of its components,  $Y_{ij}$ , has a corresponding univariate normal distribution, with conditional mean  $\mu_{ij}$  and conditional variance  $\sigma^2_{ij}$ .

Recall that while observations from different individuals are assumed to be independent of one another, repeated measurements of the same individual are not assumed to be independent. This lack of independence is captured by the off-diagonal elements of the covariance matrix,  $\Sigma_i$ . The covariance matrix has been indexed by  $i$ , and this allows, in principle, the covariance matrix to depend on the covariates,  $X_i$  (e.g., on the times of the repeated measures). In the case where all individuals have the same number of repeated measures, obtained at a common set of occasions, and where there is no dependence of the covariance matrix on the covariates, we can drop the index  $i$  and simply denote the covariance matrix by  $\Sigma$ . This would be analogous to the assumption of homogeneity of variance in linear regression for a univariate response, that is, for the vector of responses, it is assumed that there is homogeneity of covariance. However, when individuals have unequal numbers of repeated measures and/or when the repeated measures are obtained at different occasions, the covariance matrix will typically depend on the number and timing of the measurements. In principle, the covariance can also depend on covariates other than time; for example, the covariance could depend on the treatment group. However, in practice, this type of dependency of the covariance on covariates is very rarely ever assumed; this is analogous to the

ordinary univariate regression setting where we usually do not allow the error variance to depend on covariates.

# Multivariate Normal Distribution

So far we have discussed the multivariate normal distribution in a very general way, noting how it can be seen as a very natural, multivariate extension of the univariate normal distribution. Next we present a more detailed description of the multivariate normal distribution, since it forms the basis for a general method of estimation that will be described in Chapter 4. However, we remind the reader that the statistical methods presented in later chapters use the assumption that the longitudinal responses have an approximate multivariate normal distribution to derive estimates and statistical tests but do not require it when data are complete (i.e., no missing data), or when there is missingness but the observed data can be regarded as a random sample of the complete data.

The foundation for much of statistics is based on probability theory. Indeed, the formal basis for many statistical methods is an assumed probability distribution for the response variable. Broadly speaking, a probability distribution describes the likelihood or relative frequency of occurrence of particular values of the response variable. In particular, the probability density function for  $Y$ , denoted hereafter by  $f(y)$ , describes the probability or relative frequency of occurrence of particular values of  $Y$ . Before we describe some of the properties of the multivariate normal distribution, we first review the univariate normal distribution.

Consider a single univariate response from a longitudinal study at a particular occasion, say  $Y_{ij}$ . We assume that the mean of  $Y_{ij}$  is related to the covariates by the following linear regression model:

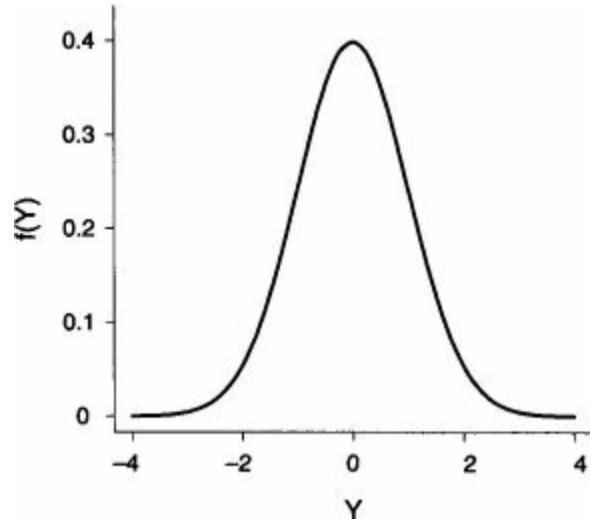
$$Y_{ij} = X'_{ij}\beta + e_{ij},$$

where the errors,  $e_{ij}$ , have a *univariate* normal distribution with mean zero and constant variance  $\sigma_j^2$ ; we denote this by  $e_{ij} \sim N(0, \sigma_j^2)$ . Recall that if the  $e_{ij}$ 's have a normal distribution with mean zero and constant variance  $\sigma_j^2$ , then  $Y_{ij}$  also has a conditional distribution that is normal, except with mean  $\mu_{ij} = X'_{ij}\beta$  and constant variance  $\sigma_j^2$ . Mathematically the univariate normal (or Gaussian) probability density function for  $Y_{ij}$  given  $X_{ij}$  can be expressed as

$$f(y_{ij}) = (2\pi\sigma_j^2)^{-1/2} \exp\left\{-\frac{1}{2}(y_{ij} - \mu_{ij})^2 / \sigma_j^2\right\},$$

where  $-\infty < y_{ij} < \infty$ . Specifically,  $f(y_{ij})$  describes the familiar bell-shaped curve illustrated in [Figure 3.1](#). Note that the area under the curve between any two values represents the probability of  $Y_{ij}$  taking a value within that range.

[\*\*Fig. 3.1\*\*](#) Plot of univariate normal density function with zero mean and unit variance.



The normal distribution has some notable features. First, the distribution is completely determined by two parameters, the mean  $\mu_{ij}$  and variance  $\sigma_j^2$  (or standard deviation  $\sigma_j$ ). Also note that the expression for the normal probability density given above depends to a very large extent on

$$\frac{(y_{ij} - \mu_{ij})^2}{\sigma_j^2} = (y_{ij} - \mu_{ij})(\sigma_j^2)^{-1}(y_{ij} - \mu_{ij}).$$

The latter is the squared distance of  $y_{ij}$  from  $\mu_{ij}$ , but expressed in standard deviation units. Thus it can be interpreted as the standardized distance of  $y_{ij}$  from its conditional mean, relative to the variability

or spread of values around the conditional mean,  $\mu_{ij}$ .

In the context of a longitudinal study, with  $n_i$  repeated measures on the  $i^{th}$  individual, we have a vector of responses and need to consider their *joint* probability distribution. While a univariate probability density function describes the probability or relative frequency of occurrence of particular values of a single random variable, a joint probability density function describes the probability or relative frequency with which the vector of responses take on a particular set of values. The multivariate normal distribution is a natural extension of the univariate normal distribution for a single response to a vector of responses. The multivariate normal joint probability density function for  $Y_i$  given  $X_i$  can be expressed as

$$f(y_i) = f(y_{i1}, y_{i2}, \dots, y_{in_i}) \\ = (2\pi)^{-n_i/2} |\Sigma_i|^{-1/2} \exp \left\{ -\frac{1}{2} (y_i - \mu_i)' \Sigma_i^{-1} (y_i - \mu_i) \right\},$$

where  $-\infty < y_{ij} < \infty$  for  $j = 1, \dots, n_i$ ,  $\mu_i = E(Y_i|X_i) = (\mu_{i1}, \dots, \mu_{in_i})'$ ,  $\Sigma_i = \text{Cov}(Y_i|X_i)$ , and  $|\Sigma_i|$  denotes the *determinant* of  $\Sigma_i$ . The determinant of  $\Sigma_i$  is also known as the *generalized variance*. The determinant of  $\Sigma_i$  summarizes the salient features of the variation expressed by  $\Sigma_i$  in a single number; a more detailed definition of  $|\Sigma_i|$  requires a greater understanding of matrix algebra than is assumed in this book.

Note the remarkable similarity between the expressions for the univariate and multivariate normal probability density functions. In some sense the multivariate normal joint probability density function simply replaces the expression for the standardized distance of  $y_{ij}$  from  $\mu_{ij}$ ,

$$(y_{ij} - \mu_{ij})(\sigma_j^2)^{-1}(y_{ij} - \mu_{ij}),$$

with a multivariate analogue for the standardized distance of the vector  $y_i$  from the vector  $\mu_i$ ,

$$(y_i - \mu_i)' \Sigma_i^{-1} (y_i - \mu_i),$$

where  $\Sigma_i^{-1}$  denotes the inverse of the matrix  $\Sigma_i$  (inversion for matrices is the analogue of the reciprocal for numbers in the sense that multiplication by the matrix  $\Sigma_i^{-1}$  can be thought of as division by the matrix  $\Sigma_i$ ). Although the latter expression is somewhat more complicated than in the univariate case, it does have interpretation in terms of a standardized measure of distance in multivariate space. For example, if  $Y_i$  is bivariate, with  $Y_i = (Y_{i1}, Y_{i2})'$ , then

$$(y_i - \mu_i)' \Sigma_i^{-1} (y_i - \mu_i) \\ = (1 - \rho_{12}^2)^{-1} \left\{ \frac{(y_{i1} - \mu_{i1})^2}{\sigma_1^2} + \frac{(y_{i2} - \mu_{i2})^2}{\sigma_2^2} - 2\rho_{12} \frac{(y_{i1} - \mu_{i1})(y_{i2} - \mu_{i2})}{\sqrt{\sigma_1^2 \sigma_2^2}} \right\},$$

where

$$\rho_{12} = \frac{\sigma_{12}}{\sqrt{\sigma_1^2 \sigma_2^2}}.$$

Although this is a more complex expression, since it accounts for the correlation between  $Y_{i1}$  and  $Y_{i2}$ , it does nonetheless provide a single measure of distance that (1) adjusts for differences in the variances of  $Y_{i1}$  and  $Y_{i2}$ , by effectively down-weighting deviations from the mean when the variance is larger, and (2) adjusts for the magnitude of the correlation (or overlapping information) between  $Y_{i1}$  and  $Y_{i2}$ . When there is no correlation between  $Y_{i1}$  and  $Y_{i2}$  (and  $\rho_{12} = 0$ ), the distance of  $y_i$  from  $\mu_i$  is simply the sum of the component standardized distances. On the other hand, when there is strong positive correlation, the distance of  $y_i$  from  $\mu_i$  also includes a component that factors in whether  $y_{i1}$  and  $y_{i2}$  are *both* larger (or smaller) than  $\mu_{i1}$  and  $\mu_{i2}$ , respectively. The latter adjustment is made because part of the standardized distance of  $y_{i2}$  from  $\mu_{i2}$  is predictable from  $Y_{i2}$ 's correlation with  $Y_{i1}$ , and vice versa.

In addition many of the properties of the multivariate normal distribution are similar to the univariate normal distribution. First, it is completely determined by the mean response vector,  $\mu_i$ , and by the covariance matrix,  $\Sigma_i$ . Also, as mentioned already,  $f(y_i)$  depends to a very large extent on the standardized distance

$$(y_i - \mu_i)' \Sigma_i^{-1} (y_i - \mu_i).$$

While the multivariate normal distribution shares many of the properties of the univariate normal distribution, the assumption of multivariate normality is much more difficult to verify from the data at hand. Unlike the univariate case, where there are simple graphical tools for assessing the validity of the assumption of normality (e.g., histograms and normal quantile plots), for most practical purposes it is difficult to assess whether a vector of responses has a conditional distribution that is multivariate normal. Although statistical tests of multivariate normality have been developed, in general, they are not very helpful because they will often detect departures from normality that are of no real substantive importance.

Perhaps the most useful assessment of the validity of the assumption of multivariate normality is through the use of graphical displays. For example, histograms and box and whisker plots of the residuals at each occasion can be used to detect gross departures of  $e_{ij}$  from univariate normality (we discuss residual diagnostics in greater detail in Chapter 10). These simple graphical displays can also be used to determine an appropriate transformation of the response variable so that the marginal distributions of the  $e_{ij}$ 's more closely approximate normal distributions. However, a caveat of the use of this technique is that although a multivariate normal distribution for  $e_i$  implies that each of the separate  $e_{ij}$  has a univariate normal distribution, univariate normal distributions for the  $e_{ij}$  do not necessarily imply that  $e_i$  has a multivariate normal distribution. Therefore the assumption of multivariate normality cannot be formally verified by examination of each of the component variables separately. However, gross departures from univariate normality can be taken to indicate that the multivariate normal assumption is not tenable.

Another property of the multivariate normal distribution for  $Y_i$  given  $X_i$  is that the association between any pair of responses is *linear*. Consequently, if the conditional distribution of  $Y_i$  is multivariate normal, then scatterplots of the residuals at all possible pairs of occasions should provide no evidence of discernible departures from a linear trend among the pairs of variables. Once again, a caveat of this simple graphical technique is that it cannot be used to establish that the conditional distribution of  $Y_i$  is multivariate normal; it can only provide evidence of discernible departures from multivariate normality.

At this point the reader may have some concerns about making the assumption of multivariate normality, especially given the inherent difficulties of verifying this assumption from the longitudinal data at hand. Fortunately, as will be discussed in later chapters, the assumption of multivariate normality is not so critical for estimation and valid inferences about  $\beta$  when data are complete (i.e., no missing data). Moreover this property extends to the setting of incomplete data if the observed data can be regarded as a random sample of the complete data. Some hints for why the normality assumption is not so critical can be found in the literature on linear regression for a single response variable. We remind the reader that there are some well-known results from linear regression models for a univariate response concerning the impact of departures from a (univariate) normal distribution. In that setting the assumption of univariate normality has been found to be not quite so critical as the assumptions made about the independence of the errors and homogeneity of the variance of the errors. That is, in linear regression for a single response it is departures from the assumption about the independence of the observations and the assumption of constant variance of the errors that have a major impact on the analysis. Departures from normality, unless they are very extreme (e.g., highly skewed response data), are not so critical. In the longitudinal data setting there are very similar results, which suggests that it is the assumptions about the dependence among the errors and assumptions about the variances and covariances that have the greatest impact on statistical inference. Departures from multivariate normality, unless they are very extreme, are not so critical. In later chapters we will discuss this topic at greater length and also describe how the assumption that  $Y_i$  has a multivariate normal distribution can be relaxed or avoided altogether.

In summary, in longitudinal studies the repeated measurements on the same individual are inherently dependent or correlated. This lack of independence can be accounted for by considering

the multivariate distribution of the entire vector of repeated measurements (given the covariates). Note that while the repeated measurements are correlated, we implicitly assume that the vectors of observations from different individuals are independent of one another. In Part II of this book we are primarily concerned with longitudinal data that are continuous, and we make the assumption that their joint distribution is multivariate normal for the purpose of deriving estimates and statistical tests. However, the methods that are discussed do not require the assumption that the responses have a multivariate normal distribution. In practice, longitudinal data are not anticipated to have a joint distribution that is *exactly* multivariate normal. The multivariate normal distribution is adopted as an approximation, but one that has many convenient statistical properties. In Chapter 4 we present a method for estimating  $\beta$  and  $\Sigma_i$ , and for making inferences about  $\beta$  and  $\Sigma_i$ , which is derived from the multivariate normal assumption for the longitudinal responses. In later chapters we discuss how the assumption that  $Y_i$  has a multivariate normal distribution can be avoided altogether.

One final comment on notation and terminology. For the remainder of the book, for simplicity of notation, we often replace  $E(Y_i|X_i)$  in many equations by  $E(Y_i)$  when it should be clear from the context that it denotes the *conditional* mean of the responses given the covariates. Likewise we often replace  $\text{Cov}(Y_i|X_i)$  by  $\text{Cov}(Y_i)$  when it should also be clear from the context that it denotes the *conditional* covariance of the responses given the covariates. In a similar vein, in discussing the distribution of  $Y_i$ , it should be understood that we are always referring to the *conditional* distribution of  $Y_i$  given the covariates.

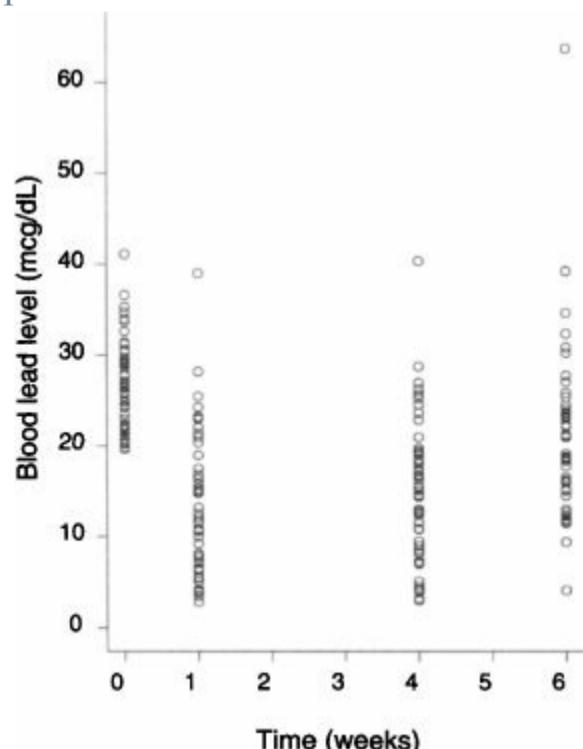
### 3.3 SIMPLE DESCRIPTIVE METHODS OF ANALYSIS

Next we consider some simple graphical tools for describing the most salient features of longitudinal data. The formal statistical analysis of longitudinal data should always be preceded by simple graphical displays of the data. A natural way to display longitudinal data is through the use of a *time plot*. A time plot is simply a scatterplot, with the responses on the vertical axis and the measurement times on the horizontal axis. For a variety of reasons the time plot of the raw longitudinal data is not always very helpful or readily interpretable. First, in most longitudinal studies the set of measurement occasions is common to many, if not all, of the study participants. As a result a time plot will result in many overlapping data points at each measurement occasion. The most extreme example of this problem arises in the time plot of binary data; it is impossible to discern any information about time trends from the resulting time plot due to the overlapping data points (e.g., 0's and 1's) at each measurement occasion.

Also note that the time plot does not indicate which data points represent repeated measurements on the same individual. To circumvent the latter problem, the time plot can be supplemented by joining or connecting successive repeated measures on the same individual with straight lines. However, the resulting line segments do not necessarily enhance the time plot; more often than not, the result is a “spaghetti” plot that is not very informative about trends in the mean response over time. Perhaps the only useful source of information provided by the simple time plot concerns the presence of extreme outliers in the data and whether the variability in the data changes discernibly with time.

Some of the problems with the time plot of longitudinal data can be illustrated using data from the *Treatment of Lead-Exposed Children Trial*. [Figure 3.2](#) displays a time plot of the blood lead level data for the group treated with succimer. Recall that in this study repeated measurements were taken at the same set of occasions for all subjects in the trial. Because of the resulting overlap of data points at the four measurement occasions, it is difficult to discern any pattern in the mean response trend over time. As noted earlier, the most extreme case of this problem arises when the response variable is binary; then it is impossible to discern any information about time trends from the resulting time plot due to the completely overlapping data points.

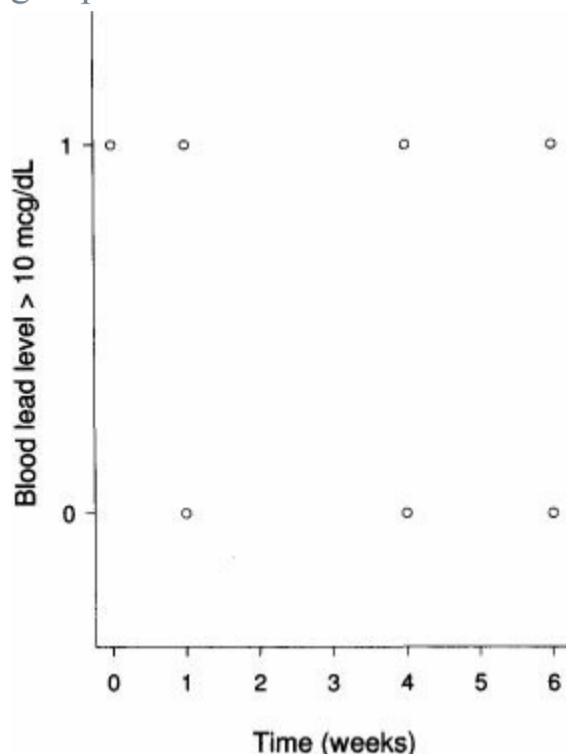
**Fig. 3.2** Time plot of blood lead levels at baseline (week 0), week 1, week 4, and week 6 for children from the succimer group in the TLC trial.



[Figure 3.3](#) displays a time plot of the repeated binary response, indicating whether each child has a blood lead level below 10  $\mu\text{g}/\text{dL}$ . (The U.S. Centers for Disease Control defines 10  $\mu\text{g}/\text{dL}$  as the

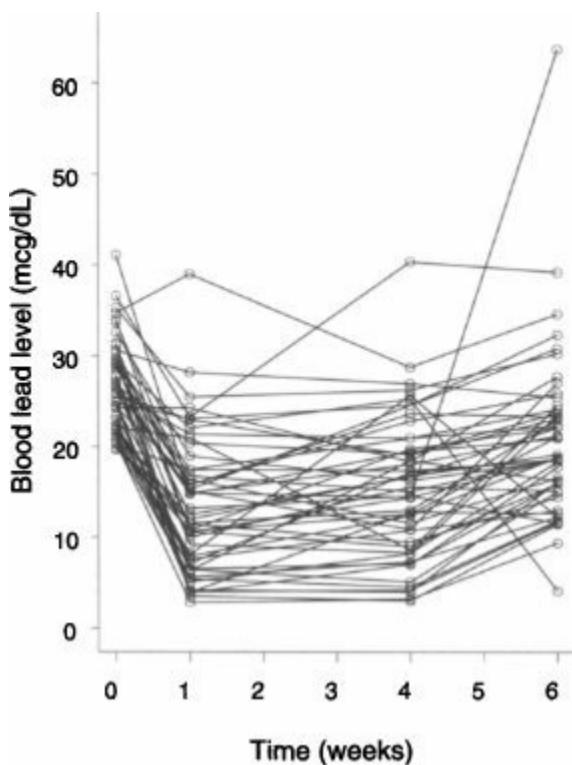
threshold for concern about exposure to lead.) Here, the binary response  $Y_{ij} = 1$  if the  $i^{th}$  child's blood lead level is above  $10 \mu\text{g}/\text{dL}$  at the  $j^{th}$  occasion, and  $Y_{ij} = 0$  otherwise. Due to the overlapping 0's and 1's, the time plot provides no information about the trend in the mean response (or probability that a blood lead level is above  $10 \mu\text{g}/\text{dL}$ ) over time.

**Fig. 3.3** Time plot of blood lead levels  $> 10 \mu\text{g}/\text{dL}$  at baseline (week 0), week 1, week 4, and week 6 for children from the succimer group in the TLC trial.



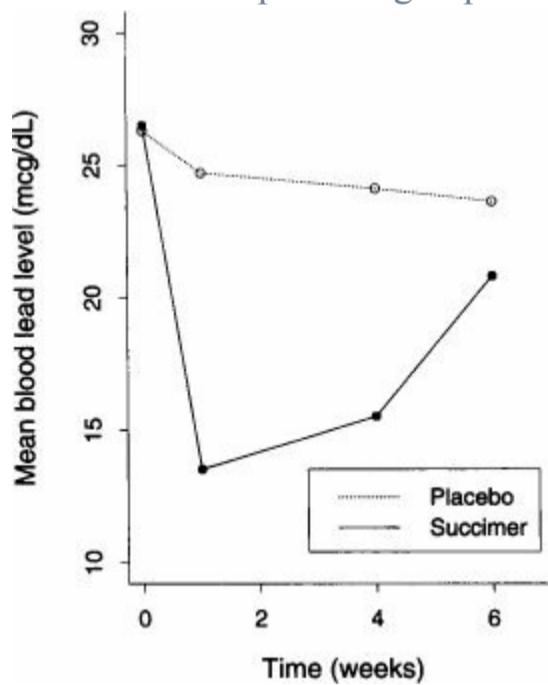
In [Figure 3.4](#) the time plot of blood lead levels is supplemented with line segments joining successive measures on the same individual. However, [Figure 3.4](#) is only a little more informative about trends in the mean response over time than [Figure 3.2](#). Although, in principle, [Figure 3.4](#) distinguishes two sources of variability in the data, between-subject variability and within-subject variability, in practice, it is difficult to assess their relative magnitude from the time plot. However, [Figure 3.4](#) does reveal an observation at week 6 that is a potential outlier, given the previous measurements of blood lead levels for that child. In summary, time plots of the raw data, with or without joined line segments for successive repeated measurements on the same individual, can reveal important features of the data. However, time plots of the raw data are not always the most informative displays of longitudinal data, especially when the data are balanced over time. With time plots of balanced longitudinal data it can be difficult to discern the “signal” (i.e., the trend in the mean response over time) from the “noise” in the data and the between-subject and within-subject sources of variability are often almost completely obscured. With highly unbalanced data, time plots of the raw data, with joined line segments, are easier to interpret.

**Fig. 3.4** Time plot, with joined line segments, of blood lead levels at baseline (week 0), week 1, week 4, and week 6 for children from the succimer group in the TLC trial.



In general, it is usually more informative to display a time plot of the average or mean response, with successive points on the graph joined by straight lines. In addition time plots of the mean response for different levels of discrete covariates (e.g., different treatment or exposure groups) can be overlayed on the same graph. The construction of these plots is relatively straightforward when the timing of the repeated measures is the same for all individuals. The time plots can also be enhanced by including standard error bars for the mean response at each occasion. For example, [Figure 3.5](#) displays the mean blood lead levels in the succimer and placebo groups at weeks 0, 1, 4, and 6. From this simple display it is readily apparent that the effect of succimer is greater after one week of treatment and that there appears to be a rebound effect thereafter.

**Fig. 3.5** Time plot, with joined line segments, of the mean blood lead levels at baseline (week 0), week 1, week 4, and week 6 in the succimer and placebo groups.



Overall, a graphical display of the mean response can be quite enlightening and can provide the basis for choosing an appropriate model for the analysis of change over time. For example, the time plot of the mean response in [Figure 3.5](#) suggests that the analysis of the blood lead levels at all four occasions may require non-linear (e.g., quadratic) or perhaps piecewise linear trends over time.

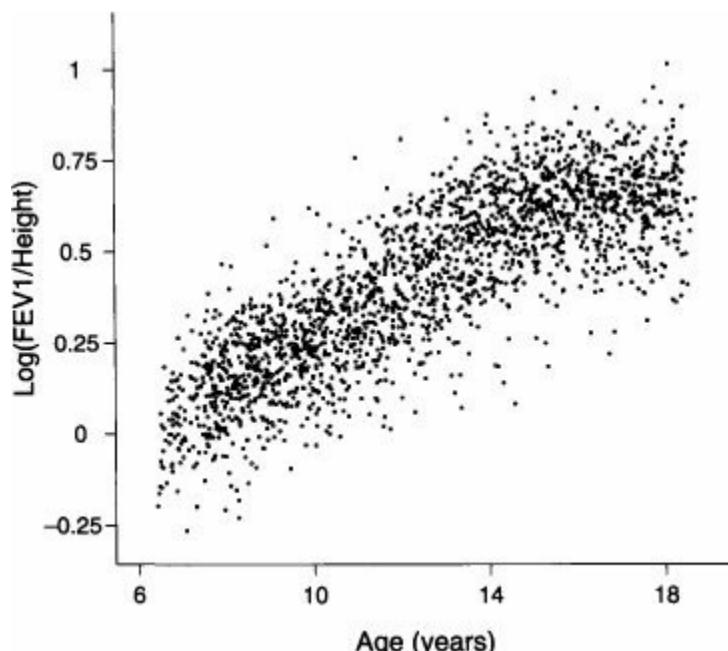
Simple time plots of the mean response are less straightforward when a covariate of interest is quantitative (e.g., dose of drug). For the purposes of producing a graphical display of the mean response trend, one simple, but often quite effective, approach is to construct a small number of groupings or “reference categories” for the quantitative covariate in question. For ease of exposition, we consider three groupings that can be generically denoted as “low,” “medium,” and “high.” Then, given this set of reference categories, the construction of the time plot of the mean response trend can

proceed along exactly the same lines as for the case of a truly discrete covariate having only three levels. That is, we can simply plot the mean response trends overlayed for the different values of the reference categories. Thus, for all practical purposes, the graphical display of the mean response trends is no more difficult when the covariate of interest is quantitative. The only question that remains is how best to choose appropriate reference categories for a quantitative covariate.

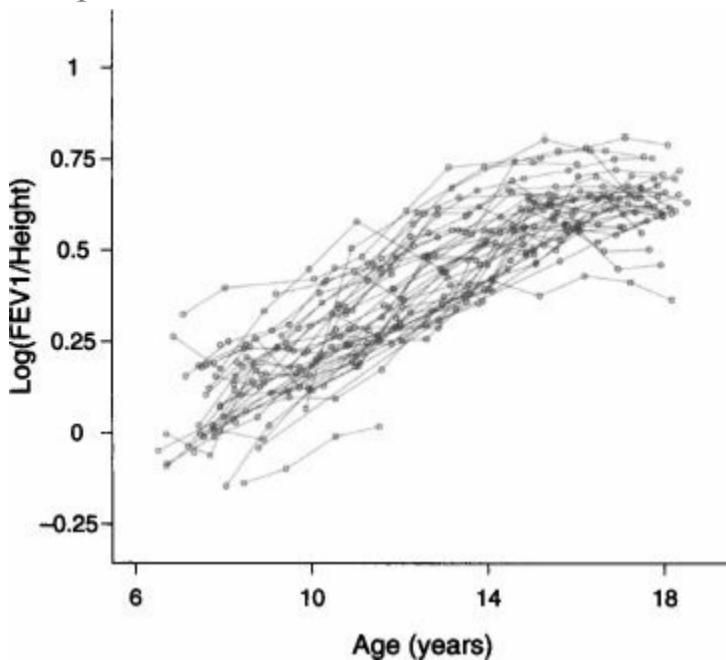
Ideally at least two or three reference categories for a quantitative covariate should be chosen, and in such a way that investigators in the field can readily appreciate the substantive importance of going from one level to the other. For example, the change in going from “low” to “medium” to “high,” or vice versa, should have some subject-matter meaning. For some quantitative covariates, there may be natural choices for the reference levels (e.g., corresponding to intervals that represent “normal” and “abnormal” ranges). For other quantitative covariates, especially those that are less well established or unfamiliar to investigators in the field, there may not be an obvious choice for the reference categories. In the latter case the choice can be made on the basis of the data at hand. For example, one possible choice is to group the covariate at the 25th and 75th percentiles. This will produce “low” (or lowest quartile), “medium” (2nd or 3rd quartiles), and “high” (or highest quartile) reference categories. It must be acknowledged, though, that the number and choices of reference groups is, to some extent, arbitrary; reference groups that are more or less extreme than those suggested here could equally be chosen (e.g., tertiles or quintiles).

So far our discussion has assumed that many, if not all, individuals are measured at the same set of occasions. When the times of measurement are not the same for all individuals, construction of time plots of the mean response can pose difficulties due to sparseness of data at any particular occasion. For example, [Figure 3.6](#) displays a time plot of longitudinal data on lung function growth in children and adolescents from the Six Cities Study of Air Pollution and Health. The data are from a cohort of 300 school-age girls living in Topeka, Kansas, who, in most cases, were enrolled in the first or second grade (between the ages of six and seven). The girls were measured annually until high school graduation (approximately at age eighteen) or loss to follow-up, and each girl provided a minimum of one and a maximum of 12 observations. At each examination, pulmonary function measurements were obtained from simple spirometry. The basic maneuver in simple spirometry is maximal inspiration followed by forced exhalation as rapidly as possible into a closed chamber. A widely used measure computed from simple spirometry is the volume of air exhaled in the first second of the maneuver, FEV<sub>1</sub>. [Figure 3.6](#) displays a time plot of log(FEV<sub>1</sub>/height) versus age for the 300 girls. Although children were measured approximately annually, the data are highly unbalanced when age, rather than chronological time, is used as the metamer for lung function growth. [Figure 3.7](#) displays a time plot, with joined line segments, of log(FEV<sub>1</sub>/height) versus age for 50 randomly selected girls. Because each girl is not measured at the same age, construction of plots of the mean response versus age can pose difficulties due to sparseness of data at any particular age.

**Fig. 3.6** Time plot of log(FEV<sub>1</sub>/height) versus age in years for girls from Topeka.



**Fig. 3.7** Time plot, with joined line segments, of  $\log(\text{FEV}_1/\text{height})$  versus age in years for 50 randomly selected girls from Topeka.



In cases where the occasions of measurement are different, it is helpful to produce a “smoothed” plot of the mean response trend over time. A smooth plot of the trend can be obtained using a variety of different approaches that can be generically referred to as “smoothing techniques.” Many of these smoothing techniques approach the estimation of the mean response at any time by considering not only the observations at that occasion but also the “neighboring” observations. That is, the estimated mean is based on observations taken before, at, and after the time of interest. The mean response at any time, say  $t$ , is taken to be a weighted average of the observations in some close proximity or neighborhood of time  $t$ .

One well-known special case of this approach is the so-called “running average” or “moving average.” For longitudinal data that are balanced and complete (no missing data), the moving average at time  $t$ , denoted  $S_t$ , is given by

$$S_t = \frac{1}{N} \sum_{i=1}^N \sum_{j=-k}^k w_j y_{i,t+j}, \quad t = k+1, \dots, n-k;$$

where  $k$  is some positive integer (e.g.,  $k = 1$  or  $k = 2$ ) and we refer to  $2k + 1$  as being the *order* of the moving average. This expression for the moving average assumes that all  $N$  individuals are measured at the same set of occasions. With highly unbalanced and/or incomplete longitudinal data, a similar expression for the moving average can be derived. The order of the moving average determines a symmetric neighborhood of values used to estimate the mean response at time  $t$ . The higher the order of the moving average, the greater is the smoothness of the resulting estimate of the mean time trend. Correspondingly, the lower the order of the moving average, the greater is the roughness of the resulting estimate of the time trend, often producing a curve that has many “wiggles” (for lack of a better term) and/or a somewhat jagged appearance. The  $w_j$  are a set of weights whose only restriction is that they must sum to one (i.e.,  $\sum_{j=-k}^k w_j = 1$ ). Ordinarily the  $w_j$  are positive, and in cases where they are unequal, they are chosen so that they decrease symmetrically about some maximum value; that is,  $w_j = w_{-j}$ , and  $w_0 > w_1 > \dots > w_k$ . As a result observations obtained in close proximity (in a temporal sense) to time  $t$  have the greatest impact or “weight” in the calculation of the mean or average response at time  $t$ . This definition of the moving average will be somewhat problematic at the beginning and end of the time plot, since the “neighborhood” of values near the end points is necessarily smaller. This problem can be rectified by altering the summation to range from  $j = \max(-k, 1-t)$  to  $j = \min(k, n-t)$  and dividing by the corresponding sum of the included weights. A simple example of a “moving average” is

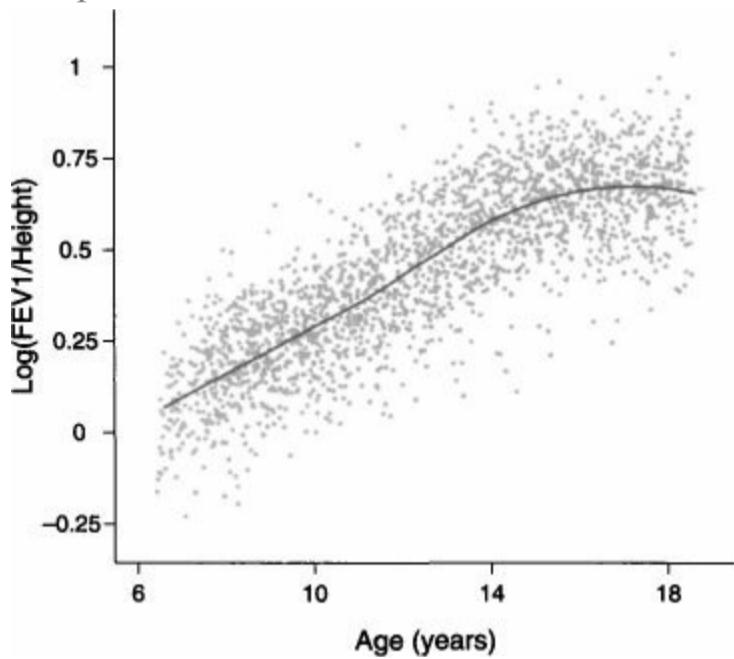
$$S_t = \frac{1}{N} \sum_{i=1}^N \frac{y_{i,t-1} + y_{it} + y_{i,t+1}}{3}.$$

In this example the weights are all equal (i.e.,  $w_{-1} = w_0 = w_1 = 1/3$ ).

Moving averages are best suited to smoothing observations that are approximately equally

separated in time. They are not ideal for handling completely irregularly spaced observations. When longitudinal data are irregularly spaced and unbalanced over time, other nonparametric regression methods can be used to estimate the mean response trend over time. One popular method available in most standard statistical software packages is locally weighted regression or *lowess*. The basic idea behind most of the nonparametric regression methods is very similar. They attempt to trace the salient features of the mean response as a function of time while making only minimal assumptions about the form of the relationship. For example, the lowess estimate at time  $t$  is best understood by imagining that there is a “window” centered at time  $t$ . One estimate of the mean at time  $t$  is obtained by taking some weighted average of all the observations that fall within the window. The lowess estimate, however, is not based on a simple weighted average of the observations within the window. Instead, it is determined by fitting a straight line to the data within the window using a robust regression technique that gives more weight to observations close to the center of the window and that also down-weights potential outliers. The lowess estimate of the mean at time  $t$  is simply the predicted value at time  $t$  from the fitted regression line. The entire lowess curve is obtained by moving a window of fixed width from the first measurement occasion to the last, and repeating the process at every time. [Figure 3.8](#) displays a lowess curve for the lung function data described earlier. Unlike the time plot of the raw data in [Figure 3.6](#), the lowess curve is informative about changes in lung function as the children grow older. The smooth curve produced by the lowess procedure indicates ages where lung function appears to develop more rapidly.

**Fig. 3.8** Time plot of  $\log(\text{FEV}_1/\text{height})$  versus age in years, with *lowess* smoothed curve superimposed, for girls from Topeka.



All smoothing techniques require that a smoothing parameter, often referred to as the *bandwidth* parameter, be specified. This parameter controls the amount of smoothing. For example, the width of the window in the lowess procedure determines how jagged or smooth the resulting plot appears; the wider the window, the smoother the resulting curve will be. The choice of smoothing or bandwidth parameter involves the classical trade-off between bias and precision. Excessive smoothing decreases the variance of the estimate of the mean trend but at the risk of introducing bias. Insufficient smoothing is unlikely to introduce bias but will result in a quite variable estimate of the mean response trend. All smoothing techniques must compromise in some way, and the goal is to find an appropriate trade-off between these two competing forces: increased bias versus decreased variance of the estimated mean response trend over time.

Finally, we note that standard applications of nonparametric smoothing techniques to longitudinal data ignore the correlation among repeated measures on the same individual. On the whole, the correlation among repeated measures is not likely to grossly distort the estimated mean response trend when standard smoothing techniques are applied to longitudinal data. Correlation is likely to have a much greater impact on the construction of confidence bands for the smoothed curve. As a result we caution the reader that the confidence bands produced by standard statistical software for lowess, and other nonparametric smoothing techniques, are likely to be optimistically biased. That

is, confidence bands constructed under the assumption that the correlation among repeated measures is zero will be too narrow and could potentially lead to misleading inferences. In summary, the routine application of nonparametric smoothing methods to longitudinal data can be useful for exposing trends in the mean response over time, with the caveat that confidence bands produced by standard statistical software packages should be ignored. A final note concerns attrition over time. If attrition or dropout occurs in a substantial number ( $> 5\%$ ) of subjects, then the end of the estimated mean curve can be distorted if subjects who leave the study differ from those who remain; the impact of dropout, and of missingness more generally, is discussed in greater detail in Section 4.3 and Chapter 17.

## 3.4 MODELING THE MEAN

In this section we introduce several approaches for modeling the mean of a vector of longitudinal responses. Two main approaches are distinguished: the analysis of response profiles and parametric or semiparametric curves. Both of these approaches are discussed in greater detail in Chapters 5 and 6.

As mentioned earlier, the analysis of longitudinal data focuses on changes in the mean response over time, and on the relation of these changes to covariates. For example, in the *Treatment of Lead-Exposed Children Trial*, the investigators were primarily interested in how blood lead levels changed over time and whether these changes were related to the treatment assigned. The fact that measurements obtained on the same individual are not independent, but are positively correlated, is an important consideration in their analysis, but for most longitudinal studies the correlation is not usually of scientific interest per se.

In Section 2.4 we mentioned that regression models for longitudinal data can usually be formulated to encapsulate the main research questions of interest in terms of a set of regression parameters. That is, certain regression parameters will have interpretations that bear directly on the scientific question or questions of interest. For example, in a regression model for the longitudinal blood lead level data from the TLC trial, treatment group interaction effects have direct interpretation in terms of how the underlying rate of change in mean blood lead levels differs between the two treatment groups. Before discussing approaches for modeling the mean response over time, it is important to clarify the distinction between *substantive* and *nuisance* parameters in the context of a longitudinal study.

# Substantive and Nuisance Parameters for Longitudinal Data

In regression models for longitudinal data, the regression parameters,  $\beta$ , relate changes in the mean response over time to covariates and are usually considered to be of primary or intrinsic interest. The regression parameters,  $\beta$ , can be defined so as to summarize important aspects of the research questions. As a result we often refer to these parameters as the *substantive* parameters. On the other hand, in many applications parameters that summarize aspects of the covariance or correlation among the repeated measures are considered to be of secondary interest. In statistics, parameters that are associated with these secondary aspects of the data are often referred to as *nuisance* parameters. Thus for the analysis of longitudinal data the correlation or covariance parameters are often thought of as nuisance parameters since there is no intrinsic interest in them.

By making this distinction between substantive and nuisance parameters, the covariance among longitudinal responses is, in a certain sense, regarded as a secondary aspect of the data (relative to the mean response over time). However, we must emphasize that this distinction does not imply that the covariance can be disregarded or simply ignored. Indeed, the covariance among repeated measures must be properly acknowledged to assure an appropriate method of analysis. The distinction between substantive and nuisance parameters has some important ramifications for the types of statistical methods that are adopted. For example, there will be a high premium attached to methods that yield valid estimates of the substantive parameters across a broad range of different assumptions about the nuisance parameters. In the context of longitudinal data, this implies that there will be a premium attached to methods that yield unbiased estimates of change in the mean response over time under a broad range of assumptions about the structure of the covariance among longitudinal responses.

Finally, we must also emphasize that the distinction between substantive and nuisance parameters should be determined only on subject-matter grounds. In the context of analyzing longitudinal data, the regression parameters,  $\beta$ , are typically the substantive parameters, since the primary focus is on characterizing changes in the mean response over time. The elements of  $\beta$  have this interpretation. The covariances among the repeated measures are nuisance parameters. However, in some settings where correlated data arise, there can be a complete reversal of roles. For example, with clustered data arising from a study of the familial aggregation of a disease-related outcome, parameters that summarize the dependence of the mean on certain risk factors, say  $\beta$ , are usually considered to be nuisance parameters, while the correlations among the responses for different family members are the substantive parameters of direct scientific interest. In family studies the goal is to determine if the presence of disease in a family member increases the risk of disease to relatives. The correlations among siblings and between parent-child pairs are of main interest because their relative magnitudes can be used to provide indirect evidence of genetic risk for the disease due to the sharing of the same gene pool. An additional example arises when investigators are interested in the heterogeneity of a treatment effect within a population; in that setting the variance of the treatment effect is of primary interest.

# Modeling the Mean Response over Time

Much of the focus in the analysis of longitudinal data is on the mean response. There are two broad approaches for modeling the mean response over time: the analysis of response profiles and parametric or semiparametric curves. The first approach allows arbitrary patterns in the mean response over time; it is related to a more traditional approach known in the statistical literature as “profile analysis.” In the analysis of response profiles, no specific time trend is assumed. Instead, the times of measurement are regarded as levels of a discrete factor. This approach to the analysis of longitudinal data is only applicable when all individuals are measured at the same set of occasions and the number of occasions is usually small. We describe the main features of the analysis of response profiles in Chapter 5.

A second approach is to assume a parametric curve (e.g., linear or quadratic trend) for the mean response over time. This approach can dramatically reduce the number of model parameters. By their very nature, parametric curves provide a very parsimonious description of trends in the mean response over time, and of covariate effects on the mean response over time. For example, a linear trend in the mean response can be characterized by a single regression parameter that has interpretation in terms of the constant rate of change in the mean response over time. In addition parametric curves describe the mean response as an explicit function of time. As a result, and in contrast to profile analysis, there is no necessity to require that all individuals in the study have the same set of measurement times, nor even the same number of repeated measurements.

Note that while the analysis of response profiles allows for an arbitrary pattern of mean responses over time, parametric curves impose an explicit structure on the mean responses. Although it will not always be possible to fit longitudinal data adequately with parametric curves, our experience with data from longitudinal studies suggests that in many cases the trends over time for the duration of the study are relatively simple (e.g., linear or quadratic trends in time). Alternatively, semiparametric curves (e.g., piecewise linear) can be adopted. A more detailed discussion of modeling the mean using parametric and semiparametric curves is presented in Chapter 6.

## 3.5 MODELING THE COVARIANCE

The defining feature of longitudinal data is that repeated responses are obtained on the same individuals over time and the resulting responses on the same individual are correlated. Although the correlation, or more generally, the covariance among the repeated responses, is not usually of intrinsic interest, it cannot simply be ignored. Instead, the covariance among repeated measures is an important aspect of the data that must be properly accounted for to yield valid inferences about the regression parameters of primary interest. Accounting for the correlation among repeated measures completes the specification of any regression model for longitudinal data and usually increases efficiency or the precision with which the regression parameters can be estimated. When an appropriate model for the covariance has been adopted, correct standard errors are obtained and valid inferences about the regression parameters can be made. In addition, when there are missing data, correct modeling of the covariance is often a requirement for obtaining valid estimates of the regression parameters. In Chapter 13 we also consider a method for analyzing longitudinal data that ignores the correlation among the repeated measures, for the purposes of estimation of the regression parameters, but makes an appropriate adjustment to the standard errors for the purposes of inference.

Three broad approaches to modeling the covariance among repeated measures can be distinguished: (1) unstructured covariance, (2) covariance pattern models, and (3) random effects covariance structures. The first is to allow any arbitrary pattern of covariance among the repeated measures. This results in what is ordinarily referred to as an “unstructured” covariance. That is, no explicit structure is assumed for the covariance among the repeated measures (other than the homogeneity of covariance across different individuals). Thus, when there are  $n$  repeated measures, the  $n$  variances at each occasion and the  $n \times (n - 1)/2$  pairwise covariances (or correlations) are estimated. Historically the unstructured covariance matrix has been the model of choice for the covariance in the analysis of response profiles. That is, the analysis of response profiles assumes arbitrary patterns for both the mean response over time (and their relation to covariates) and for the variances and covariances. This approach to modeling the covariance, however, is not limited to the analysis of response profiles and could equally be adopted when the mean response is modeled with parametric or semiparametric curves. There are two potential drawbacks with this approach. The first is that the number of covariance parameters can be quite large. If there are  $n$  measurement occasions, the  $n \times n$  covariance matrix has  $n \times (n + 1)/2$  unique parameters. Thus, in a longitudinal study with 10 measurement occasions, an unstructured covariance has 55 parameters (10 variances and 45 covariances). When the number of covariance parameters to be estimated is large relative to the sample size, then estimates are likely to be unstable. The second drawback of this approach is that it is only applicable when all individuals are measured at the same set of occasions. That is, it cannot accommodate mistimed measurements or, more generally, irregularly timed measurements.

Alternative approaches to modeling the covariance place structure on the covariance matrix. There are two main strategies. The first approach borrows ideas from the statistical literature on time series analysis. Time series data, in contrast to longitudinal data, arise from studies with a small number of replications or individuals (often only a single replication) and a large number of repeated measures. That is, in times series data  $N$ , the number of replications is small (often  $N = 1$ ) relative to the number of repeated measures,  $n$ . With longitudinal data, it is the reverse situation, with  $N$  being large relative to the number of repeated measures,  $n$ . Thus time series data consist of a small number of long sequences of repeated measurements, whereas longitudinal data consist of a large number of relatively short sequences of repeated measurements. Although times series data and longitudinal data are dissimilar in structure, and the data analytic goals are usually different, they do share one common feature: the repeated measures are correlated. Most of the models for the covariance in the time series literature incorporate at least one important aspect of longitudinal data: repeated measures taken closer together in time are expected to be more highly correlated than repeated measures further apart in time. This implies that the correlations decay as the time separation increases. Quite often the correlation among repeated measures is expressed as an explicit function of the time separation. In the latter case these models can be used with unequally spaced

observations. In addition many of the models for the variance assume *stationarity*, namely that the variance does not change as a function of time. Much of the statistical literature on the analysis of time series data has focused on parametric models that can adequately describe the covariance structure among the repeated measures with only a few parameters. These parsimonious models for the covariance can also be adopted for longitudinal data and are discussed in Chapter 7.

An alternative, and somewhat indirect, strategy for imposing structure on the covariance is through the introduction of *random effects*. Historically simple random effects models were one of the earliest approaches for analyzing repeated measures data. In the so-called univariate repeated measures ANOVA model, the correlation among repeated measurements is accounted for by the inclusion of a single individual-specific random effect. This effect can be thought of as a randomly varying intercept, representing an aggregation of all the unobserved or unmeasured factors that make some individuals “high responders” and some individuals “low responders.” The consequence of adding a single individual-specific random effect to every measurement on any given individual is that the resulting repeated measurements will be positively correlated. Thus the inclusion of random effects imposes structure on the covariance.

The univariate repeated measures ANOVA model has a very long history and has enjoyed widespread use in many fields of application; this model is discussed in greater detail in Section 3.6. Although the introduction of a single individual-specific random effect induces correlation among repeated measures, a feature of the model is that the resulting positive correlation is constant, and does not vary as a function of the time between any pair of repeated measurements. In addition the variance is constant over time. These constraints on the covariance structure are somewhat unappealing for longitudinal data. However, this problem can be easily remedied by the inclusion of more than one random effect. That is, the constraints on the covariance induced by the repeated measures ANOVA model can be relaxed by assuming that a subset of the regression parameters (e.g., intercepts and slopes) vary randomly across individuals. If the inclusion of a single individual-specific random effect induces positive correlation among repeated measures, albeit with a somewhat unappealing structure on the correlations, it should not come as a surprise that the inclusion of additional randomly varying coefficients induces patterns of correlation among the repeated measures that are somewhat less restrictive. In addition these models permit the variance to change over time in a smooth fashion. Indeed, random effects models provide both very flexible and parsimonious models for the covariance and are particularly well suited to handling longitudinal data that are irregularly timed. These models are discussed at length in Chapter 8.

## 3.6 HISTORICAL APPROACHES

We conclude this chapter with a brief survey of some of the earliest developments in methods for analyzing longitudinal and clustered data. Historically a variety of relatively simple methods have been developed for the analysis of repeated measures data. Some, but not all, of these happen to be special cases of the regression models for longitudinal data that are the focus of later chapters of this book. In this section we provide only a brief historical survey of some of these approaches, highlighting their relation to more general models, and noting some of their potential limitations. Many of the shortcomings of these methods alluded to here will be more readily apparent when the methods are viewed as special cases of the regression models considered in later chapters.

From a historical perspective, three methods for the analysis of repeated measures data can be distinguished: (1) univariate repeated measures analysis of variance (ANOVA), (2) multivariate repeated measures analysis of variance (MANOVA), and (3) methods based on summary measures. All three of these approaches have had varying degrees of popularity, and some are still in widespread use, in different areas of application. Many of these approaches are unnecessarily restrictive in their assumptions and their analytic goals. For example, ANOVA and MANOVA focus on comparing groups in terms of their mean response trend over time but provide little information about how individuals change over time. Also, as we will see later, ANOVA and MANOVA have numerous features that limit their usefulness for the analysis of longitudinal data. In contrast, the regression models that are discussed throughout the remainder of this book make more realistic assumptions and can address the major scientific questions of interest in a longitudinal study. For all the reasons that were outlined in Section 1.4, we view the regression paradigm as being the most useful, general, and versatile approach for analyzing longitudinal data arising from the health sciences.

# Repeated Measures Analysis by ANOVA

One of the earliest proposals for analyzing correlated responses was the repeated measures analysis of variance (ANOVA), sometimes referred to as the “univariate” or “mixed-model” analysis of variance. The analysis of variance paradigm was developed in the early part of the twentieth century by R. A. Fisher. Although many of the early applications of ANOVA were to designed experiments in agriculture, since then it has found widespread application in many other disciplines. In the repeated measures ANOVA model, the correlation among repeated measurements is assumed to arise from the additive contribution of an individual-specific random effect to each measurement on any given individual. Thus the model assumes the correlation between repeated measurements arises because each subject has an underlying (or latent) level of response that persists over time and influences all repeated measurements on that subject. This individual-specific effect is regarded as a random variable.

A notable feature of ANOVA models is that the response is related to a set of discrete covariates or factors. In the ANOVA paradigm the occasions of measurement are treated as an additional, within-subject, factor. Thus, if we let  $X_{ij}$  denote the vector of indicator variables for the study factors (e.g., treatment group, time, and their interaction), the repeated measures ANOVA model can be expressed as

$$Y_{ij} = X'_{ij}\beta + b_i + \epsilon_{ij},$$

where  $b_i$  is a random individual-specific effect and  $\epsilon_{ij}$  is a within-individual measurement error (it is implicitly assumed that  $X_{ij1} = 1$  for all  $i$  and all  $j$ ). Although both the  $b_i$  and  $\epsilon_{ij}$  are random, they are assumed to be independent of each other. Specifically, the  $b_i$  are assumed to have a normal distribution, with mean zero and variance,  $\text{Var}(b_i) = \sigma_b^2$ . The errors,  $\epsilon_{ij}$ , are assumed to also have a normal distribution with mean zero, but with variance,  $\text{Var}(\epsilon_{ij}) = \sigma_\epsilon^2$ .

Since both  $b_i$  and  $\epsilon_{ij}$  have mean zero, the model for the mean response, averaged over both sources of variability, is given by

$$E(Y_{ij}|X_{ij}) = \mu_{ij} = X'_{ij}\beta.$$

Thus, in the repeated measures ANOVA model, the response for the  $i^{th}$  individual is assumed to differ from the population mean,  $\mu_{ij}$ , by an individual-specific random effect,  $b_i$ , that persists throughout all measurement occasions, and a within-subject measurement error,  $\epsilon_{ij}$ . That is, the repeated measures ANOVA model distinguishes two main sources of variation in the data: between-subject variation,  $\sigma_b^2$ , and within-subject variation,  $\sigma_\epsilon^2$ . The between-subject variation acknowledges the simple fact that subjects respond differently; some are “high” responders, some are “low” responded, and some are “medium” responders. The within-subject variation acknowledges that there are random fluctuations that arise from the process of measurement, for example, due to measurement error and/or sampling variability.

Given these assumptions about the two main sources of variation, the covariance matrix of the repeated measurements has the following structure:

$$\text{Cov}(Y_i) = \begin{pmatrix} \sigma_b^2 + \sigma_\epsilon^2 & \sigma_b^2 & \sigma_b^2 & \dots & \sigma_b^2 \\ \sigma_b^2 & \sigma_b^2 + \sigma_\epsilon^2 & \sigma_b^2 & \dots & \sigma_b^2 \\ \sigma_b^2 & \sigma_b^2 & \sigma_b^2 + \sigma_\epsilon^2 & \dots & \sigma_b^2 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ \sigma_b^2 & \sigma_b^2 & \sigma_b^2 & \dots & \sigma_b^2 + \sigma_\epsilon^2 \end{pmatrix}.$$

The derivation of the variances and covariances is not important; a more detailed account will be given in Chapter 8 (see Section 8.1). What is important to note is that the variances at every occasion are equal,  $(\sigma_b^2 + \sigma_\epsilon^2)$ , as are the covariances,  $\sigma_b^2$ . Consequently, the correlation among any pair of repeated measures,

$$\text{Corr}(Y_{ij}, Y_{ik}) = \frac{\sigma_b^2}{\sigma_b^2 + \sigma_\epsilon^2},$$

is positive (by virtue of the fact that the variances,  $\sigma^2_b$  and  $\sigma^2_\varepsilon$ , must be positive) and constant, regardless of the time that has elapsed between the measurement occasions.

This particular covariance structure is also known as *compound symmetry* and has a randomization justification in certain repeated measures designs (e.g., split-plot designs). In an experiment where the within-subject factor is randomly allocated to subjects, randomization arguments can be made to show that the constant variance and constant correlation conditions hold (see Chapter 21 for a more detailed discussion of the randomization argument). Historically this provided an attractive justification for using the repeated measures analysis by ANOVA in randomized experiments. The randomization argument is simply not justifiable in the longitudinal data setting; measurement occasions cannot be randomly allocated to subjects. As a result the compound symmetry assumption for the covariance is often inappropriate for longitudinal data. That is, the constraint on the correlation among repeated measurements is somewhat unappealing for longitudinal data, where the correlations are expected to decay with increasing separation in time. Also the assumption of constant variance across time is often unrealistic. As mentioned earlier, our practical experience with many longitudinal studies has led to the empirical observation that the variances are rarely constant over time. Finally, as originally conceived, the repeated measures ANOVA model was developed for the analysis of data from designed experiments, where the repeated measures are obtained at a set of occasions common to all individuals, the covariates are discrete factors (e.g., treatment group and time), and the data are complete. Thus the repeated measures ANOVA could not be readily applied to longitudinal data that were irregularly spaced, incomplete, or when it was of interest to include quantitative covariates in the analysis.

Despite the somewhat unappealing structure imposed on the covariance, the requirement of a longitudinal design balanced on time, and the restriction to discrete covariates, the repeated measures ANOVA was nonetheless widely adopted for the analysis of longitudinal data. Perhaps one of the major reasons for its widespread use was because the ANOVA formulation led to relatively simple computational formulas that could be performed with a desk or pocket calculator (or indeed, with pen, paper, and a good deal of perseverance). Historically the repeated measures ANOVA was probably one of the few models that could realistically be fit to longitudinal data at a time when computing was in its infancy. However, with modern computing, and the widespread availability of statistical software for fitting a broader class of models for correlated data, there is little reason to analyze longitudinal data under the inherent limitations and constraints imposed by the repeated measures ANOVA model.

# Repeated Measures Analysis by MANOVA

Previously we described the repeated measures ANOVA for longitudinal data and noted in passing that it is sometimes referred to as the “univariate” or “mixed-model” analysis of variance. Analysis of variance was originally developed as a statistical model for independent observations, for example, observations on a single response variable obtained from independent subjects. By regarding the measurement occasions as levels of a within-subject factor, and by including a randomly varying individual-specific effect, the ANOVA model can be formulated in a way that allows for the possibility that repeated measures of the response obtained on the same individual are positively correlated. However, the repeated measures ANOVA is nevertheless conceptualized as a model for a single or univariate response variable.

In contrast, “multivariate” analysis of variance (MANOVA) is an extension of the analysis of variance model to handle cases where there are multiple response variables. That is, where ANOVA focuses on the analysis of a single response variable, MANOVA focuses on the analysis of a multivariate vector of response variables. In a certain sense MANOVA is a multivariate analogue of ANOVA.

Since MANOVA was originally developed for the analysis of a multivariate vector of response variables, it is worth emphasizing some of the distinctions between longitudinal responses and more general cases of multivariate responses. Recall that longitudinal data give rise to a vector of responses. Thus the responses in a longitudinal study are inherently multivariate. On the other hand, the multivariate responses arising from a longitudinal study are commensurate, being repeated measures over time of the same response variable. With longitudinal data, the repeated measures represent selected observations of the main features of some underlying continuous process that is potentially changing over time. This is in contrast to having a single measure of multiple, but substantively different or distinct, response variables (e.g., blood pressure, blood glucose, and LDL cholesterol levels). With more general multivariate data where we have single measurements of multiple, but distinct, response variables, there is no notion of an underlying continuum. Finally, with longitudinal data, the covariance among the repeated measures can be expected to have certain features or patterns; with more general multivariate data, there is rarely any indication of structure to the covariance matrix.

Thus MANOVA was developed to allow investigators to simultaneously analyze a single measure of multiple response variables (e.g., blood pressure, blood glucose, and LDL cholesterol levels), each of which is of interest in its own right. MANOVA also allowed investigators to examine linear combinations of the response variables, rather than the original variables themselves. Although MANOVA was developed for multiple, but substantively different, response variables, statisticians soon recognized that such data share a common feature with longitudinal data, namely that they are correlated. This led to the development of a very specific variant of MANOVA, known as repeated measures analysis by MANOVA (or sometimes referred to as multivariate repeated measures ANOVA), for the analysis of longitudinal data.

The repeated measures analysis by MANOVA is a special case of a more general approach known as *profile analysis*. The analysis of response profiles will be discussed in much greater detail in Chapter 5. Here we describe the basic idea underlying the repeated measures analysis by MANOVA, but without much technical detail. In Chapter 5 we will illustrate how the repeated measures analysis by MANOVA relates to profile analysis and highlight the potential limitations of this approach for analyzing longitudinal data.

The main idea underlying the repeated measures analysis by MANOVA can be best understood by considering a simple example. Suppose that we have two treatment groups (e.g., placebo and active treatment) and subjects are measured repeatedly on  $n$  occasions. In such a study design, three fundamental questions can be considered:

1. Are the trends in the mean response over time the same in the two groups?
2. Averaged over the two groups, is the overall trend in the mean response over time flat?

### 3. Are the overall mean responses, averaged over occasions, the same in the two groups?

Note that the first question is equivalent to asking whether there is a “group  $\times$  time interaction.” Ordinarily this first question must be addressed before consideration of the remaining questions, since it rarely makes sense to examine group or time main effects when there is an interaction. The second and third questions are equivalent to asking whether there are main effects of “time” and “group,” respectively.

To address each of these questions the repeated measures analysis by MANOVA proceeds by constructing a new set of variables, derived from the original set of repeated measures. The new set of derived variables, numbering as many as the number of repeated measures, then form the basis of a MANOVA. That is, the repeated measures analysis by MANOVA proceeds by constructing a set of derived variables and uses relevant subsets of these to address each of the three questions posed above.

A simple example will help to motivate the main ideas. Suppose that in a longitudinal clinical trial, designed to compare a new treatment to placebo, repeated measures of the response variable are obtained on three occasions ( $n = 3$ ). A repeated measures analysis by MANOVA proceeds by constructing three derived variables, say  $V_{i1}$ ,  $V_{i2}$ , and  $V_{i3}$ . The first derived variable is simply the sum (or average) of the responses. That is, for each individual we can construct

$$V_{i1} = (Y_{i1} + Y_{i2} + Y_{i3}).$$

This derived variable provides no information about within-individual changes in the response over time. Instead, it provides information about the mean level of the response, averaged over all three occasions.

The two remaining derived variables are constructed to provide information about possible within-individual changes in the response over time. For example, the following two derived variables

$$V_{i2} = (Y_{i2} - Y_{i1}) \text{ and } V_{i3} = (Y_{i3} - Y_{i1})$$

provide information about changes in the response (from time 1) at times 2 and 3, respectively. The set of three derived variables can be obtained by applying the transformation matrix

$$\begin{pmatrix} 1 & 1 & 1 \\ -1 & 1 & 0 \\ -1 & 0 & 1 \end{pmatrix},$$

to the original vector of responses. That is,

$$\begin{pmatrix} V_{i1} \\ V_{i2} \\ V_{i3} \end{pmatrix} = \begin{pmatrix} 1 & 1 & 1 \\ -1 & 1 & 0 \\ -1 & 0 & 1 \end{pmatrix} \begin{pmatrix} Y_{i1} \\ Y_{i2} \\ Y_{i3} \end{pmatrix}.$$

Thus, in the repeated measures analysis by MANOVA, a transformation matrix takes a sequence of repeated measures and produces an equal number of derived variables that are used in subsequent analyses. The first row of the transformation matrix creates the sum (or average) of the repeated measures (it makes no difference whether the sum or average is used since the latter is proportional to the former). The first derived variable provides information about the mean level of the response, averaged over all measurement occasions, and can be used to address the third question concerning whether there is a “group” effect. The remaining rows of the transformation matrix construct derived variables that provide information about change over time. There are many different ways to obtain a set of derived variables that describe change over time, and so there are many possible choices of values for the remaining rows of the transformation matrix. For example, if it is of interest to construct derived variables that represent linear and quadratic contrasts of time, the following transformation matrix can be used:

$$\begin{pmatrix} 1 & 1 & 1 \\ -1 & 0 & 1 \\ 1 & -2 & 1 \end{pmatrix}.$$

It can be shown that the multivariate statistics for tests of “time” effects and their interactions produced by the repeated measures analysis by MANOVA are invariant to how change over time is characterized in the transformation matrix.

Given the set of derived variables, the repeated measures analysis by MANOVA proceeds as follows. First, the two derived variables representing contrasts of time,  $V_{i2}$  and  $V_{i3}$ , are analyzed using MANOVA. For example, the first question can be addressed by comparing the groups in terms of these two derived variables. Specifically, in our simple example with two groups, this is achieved by using a multivariate extension of the two-sample  $t$ -test, which is known as Hotelling's  $T^2$  test. A test of no differences between groups on these two derived variables is equivalent to a test of the "group  $\times$  time interaction." Next, and assuming that there is no "group  $\times$  time interaction," the second question can be addressed. The second question is concerned with the shape of the overall (i.e., averaged over groups) trend in the mean response over time. If the mean response trend over time is flat, then the two derived variables have expectation zero. As a result the second question can be addressed by using a multivariate extension of the single sample  $t$ -test, the single sample Hotelling's  $T^2$  test, of the hypothesis that  $V_{i2}$  and  $V_{i3}$  have mean zero. Finally, the third question can be addressed by focusing on the first derived variable,  $V_{i1}$ . This variable is proportional to the mean of the repeated measures and can be used to assess whether there is a "group" effect. Specifically, group difference in the overall mean response, averaged over measurement occasions (or time), can be examined using a simple two-sample  $t$ -test, or, in the case of more than two groups, using ANOVA. A standard ANOVA of the first derived variable (the sum or average of the repeated measures) is performed and provides a test of the group or between-subject factor. Note that there is nothing intrinsically multivariate in this last part of the analysis since the analysis is based on a single derived variable. It is the first part of the analysis that is intrinsically multivariate.

In summary, the underlying idea behind repeated measures analysis by MANOVA is to obtain a new set of derived variables, based on a linear combination of the original sequence of repeated measures. The derived variables can be partitioned into a set that provides information about change, and a single derived variable that provides information about overall level of response. The latter can be analyzed using univariate ANOVA, and the results of this analysis determine whether there are "group" or between-subject effects. This analysis addresses the question of whether the groups differ in their overall level of response. The remaining derived variables are analyzed using MANOVA, and these analyses determine whether there are time effects and group  $\times$  time interactions. A test of the group  $\times$  time interactions addresses the question of whether the changes in the mean response over time are different in the groups. If there are no differences between groups in these derived variables, it is then of interest to ask whether the combined average of the derived variables, where the averaging is over groups, is different from zero. This addresses the question of whether there is any overall change in the mean response over time. In effect, this can be considered a test of the main effect of the time factor.

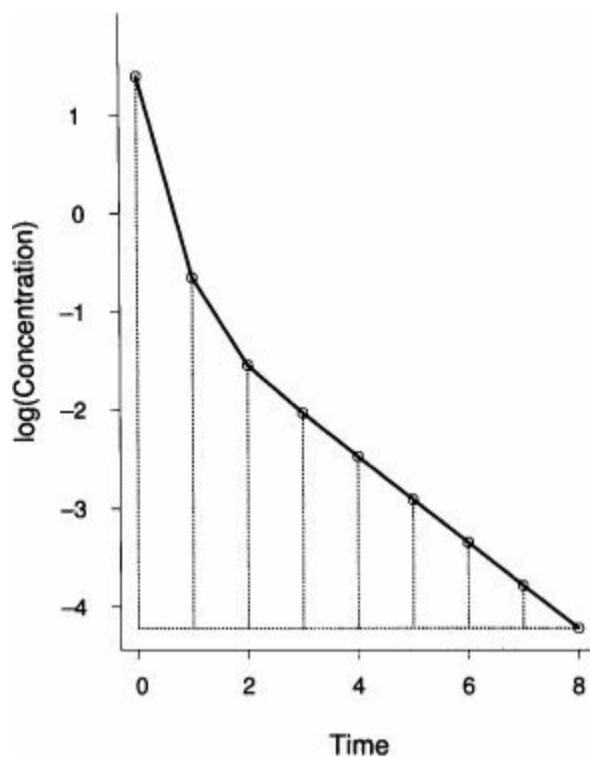
The repeated measures analysis by MANOVA has a number of features that make it unappealing for the analysis of longitudinal data. In particular, the MANOVA formulation forces the within-subject covariates to be the same for all individuals in the study. There are at least two practical consequences of this constraint. First, repeated measures MANOVA cannot be used when the design is unbalanced over time, that is, when the vectors of repeated measures are of different lengths and/or obtained at different sequences of time. Second, the repeated measures MANOVA (as implemented in existing statistical software packages) does not allow missing data. If an individual has a single missing response at any occasion, the entire data vector from that individual is excluded from the analysis. This "listwise" deletion of missing data from the analysis can result in dramatically reduced sample size and very inefficient use of the available data. Listwise deletion of missing data can also produce biased estimates of change in the mean response over time when the "completers" (i.e., those with no missing data) are not a random sample from the target population. When the "completers" are a biased sample from the target population, the sample means, variances, and covariances are biased estimates of the corresponding parameters in the target population. Some additional drawbacks of the repeated measures analysis by MANOVA will be discussed in Chapter 5, where a more detailed exposition on the analysis of response profiles is given.

# Summary Measure Analysis

A common approach to the analysis of longitudinal data still in widespread use reduces the sequence of repeated measures for each individual to a small set of summary values. The major motivation behind this approach is that if the sequence of repeated measures can be reduced to a single number summary, then standard parametric or nonparametric methods for the analysis of a univariate response can be applied to the derived measures.

For example, the area under the curve (AUC) is one common measure that is often used to summarize the sequence of repeated measures on any individual. The use of AUC is appropriate when the repeated measures for each individual are obtained at the same set of occasions. The AUC can be especially appealing in pharmacological studies where the response, or some transformation of the response, measures the absorption, concentration, or clearance of drugs. For example, the AUC can be used to estimate the clearance rate or plasma concentration of a particular dose of a drug or substance over time. The AUC can be approximated for each individual by joining adjacent measurements by line segments and summarizing the area under the curve by the sum of the areas of the resulting trapezoids (see [Figure 3.9](#)). The resulting AUC's can then be related to covariates (e.g., treatment or exposure group) using standard methods for the analysis of a univariate response (e.g., *t*-test, ANOVA, Wilcoxon rank sum test, or Kruskal–Wallis test).

**Fig. 3.9** Time plot of log(concentration) versus time, illustrating how the area under the curve (AUC) can be calculated using the trapezoidal rule.



When the covariates are discrete, and the repeated measures for each individual are obtained at the same set of occasions, the AUC analysis can also be based on the results of an analysis of response profiles (that assumes arbitrary patterns for the mean responses over time). Given the covariates, the AUC for the mean response over time is the same as the average (or mean) of the individual-specific AUCs. That is, for the case of linear models for continuous longitudinal responses, the AUC for the mean response over time coincides with the mean of the AUCs for the individuals in the population of interest. In a limited context, reducing the sequence of repeated measures for each individual to an AUC can provide a useful basis for the analysis of longitudinal data. However, the analysis of AUCs is problematic with unbalanced longitudinal data.

Another measure commonly used to summarize the sequence of repeated measures is the slope or constant rate of change in the response over time. For example, it might be assumed that a straight line (the simplest possible curve) fits the observed responses for each subject. If  $Y_{ij}$  is the response of the  $i^{th}$  individual measured at time  $t_{ij}$ , it might be assumed that

$$Y_{ij} = b_{1i} + b_{2i} t_{ij} + \epsilon_{ij},$$

where  $b_{1i}$  and  $b_{2i}$  are regression parameters specific to the  $i^{th}$  individual and the errors,  $\epsilon_{ij}$ , are

implicitly assumed to be independent within any individual. Estimates of the individual-specific slopes (and intercepts) can then be obtained from a linear regression line fit to each individual's repeated measures. The resulting slopes can then be related to the covariates using standard parametric or nonparametric methods for the analysis of a univariate response. This approach does not require that the repeated measures for each individual be obtained at the same set of occasions. Finally, we note that extensions of this particular summary measure analysis approach lead naturally to a class of models referred to as "growth curve models". Studies of growth and aging are classic examples of observational longitudinal studies. In these studies the goal is to describe naturally occurring changes in the response over time, due to developmental or aging processes, and to compare these growth curve profiles in different groups (e.g., males and females). To meet this goal, growth curve models have been developed to summarize the pattern of response over time and allow for the possibility that individuals may belong to or be drawn from different groups. Growth curve models can be motivated in terms of a two-stage model. Indeed, growth curve models are sometimes referred to as "two-stage" growth models. At the first stage, it is assumed that a parametric curve (e.g., linear or quadratic trend in time) fits the observed responses for each subject. In the second stage, these individual-specific growth parameters are then related to covariates that describe the different groups from which the individuals have been drawn. Growth curve models will be discussed in greater detail in Chapter 8.

Before the advent of modern computing and readily available statistical software for analyzing correlated data, summary measure analysis of longitudinal data had some very obvious appeal. First, the summary measures and their subsequent analysis can readily be understood by investigators with limited training in statistics. Also, once a summary measure has been derived, standard methods for the analysis of a univariate response (e.g., *t*-test, ANOVA, linear regression, Wilcoxon rank sum test, Kruskal–Wallis test) can be validly applied since issues of correlation among the observations no longer arise. That is, the summary measures on different individuals are independent of one another. Summary measure analysis can also be appealing when sample sizes are not sufficiently large for estimation of the correlation among the repeated measures. However, despite the simplicity of the method, it does have a number of distinct drawbacks. One drawback is that it forces the data analyst to focus on only a single aspect of the repeated measures over time. It should be intuitively clear that when  $n$  repeated measures are replaced by a single number summary, there must necessarily be some loss of information. Furthermore individuals with discernibly different response profiles can have the same summary measure. For example, individual-specific response profiles with quite distinct shapes can result in the same AUC. Another potential drawback of the summary measure approach is that the covariates must be time-invariant (sometimes referred to as "time-stationary" covariates). Thus, if one of the key covariates is time-varying, the method cannot be applied. Finally, we note that some of the summary measures that have been proposed are not well defined when there are missing data or irregularly spaced repeated measures. Even when they can be defined, these simple methods lose efficiency.

In those cases where the summary measures are well defined when individuals have missing data or different numbers of repeated measures, the analysis becomes more complicated because the derived summary measures no longer have the same variance. Similarly, if the repeated measures are taken at irregular times for different individuals, the resulting summary measures may also have different variances. In all these cases the variance of the derived summary measures is not constant, violating a fundamental assumption made by many standard statistical methods for univariate responses. Thus, in general, the standard parametric methods for the analysis of a univariate response (e.g., *t*-tests, ANOVA, linear regression) cannot be validly applied to the summary measures when the design is unbalanced over time due to missing data, different numbers of repeated measures, or sequences of repeated measures taken at irregular times for different individuals.

When longitudinal data are unbalanced over time, a proper analysis of the summary measures would require that each summary measure be weighted differently. However, the chief complication here is that the specific weights given to each summary measure will, in general, depend implicitly on the covariance among the repeated measures. Thus a simple univariate analysis cannot proceed

without proper consideration of the covariance, the very feature of the data that these methods were developed to avoid having to specify. In conclusion, in limited contexts summary measure analysis of longitudinal data can be useful, but it should be avoided when the data are unbalanced. When it is desirable to base analysis on a single aspect of the repeated measures over time, the regression models that are the focus of later chapters can be used. The regression modeling approach is more efficient than the summary measure analysis and can also handle unbalanced data.

## **3.7 FURTHER READING**

Winer (1971) provides a very accessible discussion of the application of repeated measures analysis by ANOVA; also see McCulloch (2005) for a comparison of repeated measures ANOVA with more modern methods of analyses. A comprehensive description of repeated measures analysis by MANOVA, targeted at applied researchers, can be found in the book by Hand and Taylor (1987). Finally, for a non-technical discussion of the analysis of summary measures, readers are referred to the review articles by Matthews et al. (1990) and Everitt (1995).

# Bibliographic Notes

An excellent discussion of scatterplot smoothing techniques can be found in Chapter 3 of Ruppert et al. (2003). Ware and Liang (1996) provide an interesting historical perspective on the development of statistical methods for the analysis of longitudinal data, with emphasis on the contributions that have been made in the biostatistical literature; also, see Chapter 1 of Fitzmaurice et al. (2009).

The foundations for the repeated measures analysis of variance can be found in the seminal monograph by Fisher (1925) and in the method for analyzing split-plot experiments proposed by Yates (1935); also see Scheffé (1959). Greenhouse and Geisser (1959) described an adjustment to the repeated measures analysis by ANOVA when the required assumption about the covariance matrix (compound symmetry) does not hold. The repeated measures analysis by MANOVA was introduced in the statistical literature by Box (1950); also see Danford et al. (1960), Geisser (1963), Cole and Grizzle (1966), and Morrison (1972). A discussion of repeated measures analyses by ANOVA and MANOVA, and the relationship between the two methods, can be found in Chapters 2, 3, and 11 of Hand and Crowder (1996).

Finally, the analysis of summary measures has a long history, dating back to the early contributions to growth curve analysis by Wishart (1938), Box (1950), and Rao (1958). Rowell and Walters (1976), in a classic paper on the analysis of longitudinal agricultural experiments, describe how linear regressions can be fitted to longitudinal data on each subject, followed by an analysis of the values of the resulting regression coefficients. The article by Rowell and Walters (1976) is widely cited for popularizing the analysis of summary measures of growth in many different disciplines.

<sup>1</sup> In all data sets used throughout this book, the original subject IDs have been replaced with new subject ID numbers to ensure that the data sets cannot be linked to the original records.

# *Chapter 4*

## *Estimation and Statistical Inference*

### **4.1 INTRODUCTION**

So far our discussion of models for longitudinal data has been very general, with no mention of methods for estimating the regression coefficients or the covariance among the repeated measures. In Chapters 5 through 8 we will consider models for longitudinal data where the response variable is continuous and assumed to have a conditional distribution that is approximately multivariate normal. In these chapters the main focus is on various aspects of modeling longitudinal data, with particular emphasis on models for the mean and covariance. All of the models presented in Chapters 5 through 8 can be expressed in terms of a general linear regression model for the mean response vector

$$(4.1) \quad E(Y_i|X_i) = X_i\beta,$$

where the response vector,  $Y_i$ , is assumed to have a conditional distribution that is multivariate normal with covariance matrix

$$(4.2) \quad \text{Cov}(Y_i|X_i) = \Sigma_i = \Sigma_i(\theta),$$

where  $\theta$  is a  $q \times 1$  vector of covariance parameters. For example, with balanced longitudinal data ( $n_i = n$ ), where an “unstructured” covariance matrix has been assumed, the elements of  $\theta$  are simply the  $n$  variances and  $\frac{n(n-1)}{2}$  pairwise covariances stacked in a single  $q \times 1$  vector (where  $\frac{n(n+1)}{2}$ ). On the other hand, if the covariance is assumed to have a “compound symmetry” pattern, then  $q = 2$  and the two elements of  $\theta$  represent the common value of the variances and common value of the pairwise covariances. In this section we consider a framework for estimation of the unknown parameters,  $\beta$  and  $\theta$  (or equivalently,  $\Sigma_i$ ).

## 4.2 ESTIMATION: MAXIMUM LIKELIHOOD

Given that full distributional assumptions have been made about the vector of responses,  $Y_i$ , since the multivariate normal distribution is entirely specified by the mean vector and covariance matrix, a very general approach to estimation is the method of *maximum likelihood* (ML). The fundamental idea behind ML estimation is really quite simple and is conveyed by its name: use as estimates of  $\beta$  and  $\theta$  the values that are most probable (or most “likely”) for the data that have actually been observed. The maximum likelihood estimates of  $\beta$  and  $\theta$  are those values of  $\beta$  and  $\theta$  that maximize the joint probability of the response variables evaluated at their observed values. The probability of the response variables evaluated at the fixed set of observed values, and regarded as functions of  $\beta$  and  $\Sigma_i(\theta)$ , is known as the *likelihood function*. Thus estimation of  $\beta$  and  $\theta$  proceeds by maximizing the likelihood function. In a certain sense the method of maximum likelihood chooses values of  $\beta$  and  $\theta$  that best explain the observed data. The values of  $\beta$  and  $\theta$  that maximize the likelihood function are called the *maximum likelihood estimates* of  $\beta$  and  $\Sigma_i(\theta)$ , and are usually denoted  $\hat{\beta}$  and  $\hat{\Sigma}_i$  (or  $\Sigma_i(\hat{\theta})$ ).

Before we present any more details concerning maximum likelihood estimation of  $\beta$  and  $\theta$ , it will be informative to consider this method of estimation in the simpler case where all observations can be assumed to be independent, that is, in the standard linear regression model with independent (and hence uncorrelated) errors that are assumed to have a univariate normal distribution.

# Independent Observations

Suppose that the data arise from a series of cross-sectional studies that are repeated at  $n$  different occasions. At each occasion, data are obtained on a sample of  $N$  individuals. Here it is reasonable to assume that the observations are independent of one another, since each individual is measured at only one occasion. Also, for ease of exposition, we assume that the variance is constant, say  $\sigma^2$ . The mean response is related to the covariates via the following linear regression model:

$$E(Y_{ij}|X_{ij}) = X'_{ij}\beta.$$

To obtain maximum likelihood estimates of  $\beta$ , we must find the values of the regression parameters that maximize the joint normal probability density function of all the observations, evaluated at the observed values of the response, and regarded as a function of  $\beta$  (and  $\sigma^2$ ). Recall from Section 3.2 that the univariate normal (or Gaussian) probability density function for  $Y_{ij}$  given  $X_{ij}$  can be expressed as

$$f(y_{ij}) = (2\pi\sigma^2)^{-1/2} \exp\left\{-\frac{1}{2}(y_{ij} - \mu_{ij})^2/\sigma^2\right\},$$

where  $-\infty < y_{ij} < \infty$ . When all the responses are independent of one another, the likelihood function is simply the product of the individual univariate normal probability density functions for  $Y_{ij}$  given  $X_{ij}$ ,

$$\prod_{i=1}^N \prod_{j=1}^n f(y_{ij}).$$

It is more common to work with the log-likelihood function, which will involve sums, rather than products, of the individual univariate normal probability density functions for  $Y_{ij}$ . Note that maximizing the likelihood is equivalent to maximizing the logarithm of the likelihood; the latter is denoted by  $l$ . Hence the goal is to maximize

$$l = \log \left\{ \prod_{i=1}^N \prod_{j=1}^n f(y_{ij}) \right\} = -\frac{K}{2} \log(2\pi\sigma^2) - \frac{1}{2} \sum_{i=1}^N \sum_{j=1}^n (y_{ij} - X'_{ij}\beta)^2 / \sigma^2,$$

evaluated at the observed numerical values of the data, with respect to the regression parameters,  $\beta$ . Here  $K = n \times N$ , the total number of observations. Note that  $\beta$  does not appear in the first term in the log-likelihood; as a result this term can be ignored when maximizing the log-likelihood with respect to  $\beta$ . Furthermore, since the second term has a negative sign, maximizing the log-likelihood with respect to  $\beta$  is equivalent to minimizing the following function:

$$\sum_{i=1}^N \sum_{j=1}^n (y_{ij} - X'_{ij}\beta)^2.$$

Maximizing or minimizing a function is a common mathematical problem that can be solved using calculus. Specifically, the maximum likelihood estimate of  $\beta$  can be obtained by equating the derivative of the log-likelihood, often called the score function, to zero and finding the solution to the resulting equation. However, in the example considered here, there is no real need to resort to calculus. Obtaining the maximum likelihood estimate of  $\beta$  is equivalent to finding the ordinary least squares (OLS) estimate of  $\beta$ , that is, the value of  $\beta$  that minimizes the sum of the squares of the residuals. Using vector notation, the least squares solution can be written as

$$\hat{\beta} = \left\{ \sum_{i=1}^N \sum_{j=1}^n (X_{ij} X'_{ij}) \right\}^{-1} \sum_{i=1}^N \sum_{j=1}^n (X_{ij} y_{ij}).$$

This least squares estimate is the value produced by any standard statistical software for linear regression (e.g., PROC GLM or PROC REG in SAS, the `lm` function in R and S-Plus, and the `regress` command in Stata). In the next section we consider how these ideas can be extended to the setting of correlated data. Also the alert reader might have noticed that we have thus far only focused on estimation of  $\beta$ , ignoring estimation of  $\sigma^2$ ; in the next section we also consider estimation of the covariance matrix.

# Correlated Observations

When there are  $n_i$  repeated measures on the same individual, it cannot be assumed that these repeated measures are independent. As a result we need to consider the joint probability density function for the vector of repeated measures. Note, however, that the vectors of repeated measures are assumed to be independent of one another. Thus the log-likelihood function,  $l$ , can be expressed as a sum of the individual multivariate normal probability density functions for  $Y_i$  given  $X_i$ .

To find the maximum likelihood estimate of  $\beta$  in the repeated measures setting, we first assume that  $\Sigma_i$  (or  $\theta$ ) is *known* (and therefore does not need to be estimated); later we will relax this very unrealistic assumption. To obtain the maximum likelihood estimate of  $\beta$ , we must find the value of  $\beta$  that maximizes the log-likelihood function. Given that  $Y_i = (Y_{i1}, Y_{i2}, \dots, Y_{in_i})'$  is assumed to have a conditional distribution that is multivariate normal, we must maximize the following log-likelihood function:

$$l = -\frac{K}{2} \log(2\pi) - \frac{1}{2} \sum_{i=1}^N \log |\Sigma_i| - \frac{1}{2} \left\{ \sum_{i=1}^N (y_i - X_i\beta)' \Sigma_i^{-1} (y_i - X_i\beta) \right\},$$

where  $K = (\sum_{i=1}^N n_i)$  is the total number of observations. Note that  $\beta$  does not appear in the first two terms in the log-likelihood; as a result these two terms can be ignored when maximizing the log-likelihood with respect to  $\beta$ . Furthermore, since the third term has a negative sign, maximizing the log-likelihood with respect to  $\beta$  is equivalent to minimizing

$$(4.3) \quad \sum_{i=1}^N (y_i - X_i\beta)' \Sigma_i^{-1} (y_i - X_i\beta).$$

The estimator of  $\beta$  that minimizes this expression is known as the *generalized least squares* (GLS) estimator of  $\beta$  and can be expressed as

$$(4.4) \quad \hat{\beta} = \left\{ \sum_{i=1}^N (X_i' \Sigma_i^{-1} X_i) \right\}^{-1} \sum_{i=1}^N (X_i' \Sigma_i^{-1} y_i).$$

Recall that so far we have made the somewhat unrealistic assumption that  $\Sigma_i$ , or  $\theta$ , is *known*. Before considering how to proceed when we must relax this assumption, it is worth discussing some of the properties of the GLS estimator of  $\beta$  when  $\Sigma_i$  is known. The first very notable property is that for any choice of  $\Sigma_i$ , the GLS estimate of  $\beta$  is unbiased; that is,

$$E(\hat{\beta}) = \beta.$$

In addition, in large samples (or asymptotically), the sampling distribution of  $\hat{\beta}$  can be shown to have a multivariate normal distribution with mean,  $\beta$ , and covariance,

$$(4.5) \quad \text{Cov}(\hat{\beta}) = \left\{ \sum_{i=1}^N (X_i' \Sigma_i^{-1} X_i) \right\}^{-1}.$$

This is true exactly when  $Y_i$  has a conditional distribution that is multivariate normal, and true in large samples even when the conditional distribution of  $Y_i$  is not multivariate normal. (By “large samples” we mean that the sample size,  $N$ , grows larger while the number of repeated measures and model parameters remains fixed.) Thus an important property of the GLS estimator of  $\beta$ , derived under the assumption of a multivariate normal distribution for  $Y_i$  given  $X_i$ , is that it provides a valid estimate of  $\beta$  even when the multivariate normal distribution assumption does not hold. Also note that if  $\Sigma_i$  is assumed to be a diagonal matrix, with constant variance  $\sigma^2$  along the diagonal (i.e., the correlations are zero and the variances are constant), the GLS estimator reduces to the ordinary least squares (OLS) estimator considered earlier. Finally, although the GLS estimator of  $\beta$  is unbiased for any choice of  $\Sigma_i$ , it can be shown that the most efficient GLS estimator of  $\beta$  (i.e., the estimator having smallest variance or greatest precision) is the one that uses the true value of  $\Sigma_i$ .

Before the reader becomes exasperated, we must address the nagging concern that we usually do not know  $\Sigma_i$  (or  $\theta$ ). Instead, we typically must estimate  $\Sigma_i(\theta)$  from the data at hand. Maximum likelihood estimation of  $\theta$  proceeds in the same way as with estimation of  $\beta$ . That is, the maximum

likelihood estimate of  $\theta$  is obtained by maximizing the log-likelihood with respect to  $\theta$ . As mentioned earlier, the problem of maximizing a function is a common mathematical problem that can be solved using calculus. Specifically, the maximum likelihood estimate of  $\theta$  can be obtained by equating the derivative of the log-likelihood with respect to  $\theta$ , also known as the score function, to zero and finding the solution to the resulting equation. However, in general, this equation is non-linear, and it is not possible to write down simple, closed-form expressions for the ML estimator of  $\theta$ . Instead, the ML estimate must be found by solving these equations using an iterative technique. Fortunately, computer algorithms have been developed to find the solution. Once the ML estimate of  $\theta$  has been obtained, we then simply substitute the estimate of  $\Sigma_i(\theta)$ , say  $\hat{\Sigma}_i = \Sigma_i(\hat{\theta})$ , into the generalized least squares estimator of  $\beta$  given by (4.4) to obtain the maximum likelihood (ML) estimate of  $\beta$ :

$$(4.6) \quad \hat{\beta} = \left\{ \sum_{i=1}^N (X_i' \hat{\Sigma}_i^{-1} X_i) \right\}^{-1} \sum_{i=1}^N (X_i' \hat{\Sigma}_i^{-1} y_i).$$

Interestingly, in large samples (or asymptotically), the resulting estimator of  $\beta$  that substitutes the ML estimate of  $\Sigma_i$  has all of the same properties as when  $\Sigma_i$  is actually known (the case we first considered). That is, in large samples:

1.  $\hat{\beta}$  is a consistent estimator of  $\beta$ ; this property can be loosely interpreted to mean that there is very high probability that  $\hat{\beta}$  is close to the population regression parameters  $\beta$  for increasing sample size  $N$ . If the distribution of the errors,  $e_i$ , is assumed to be normal, or even under the weaker assumption that the distribution of  $e_i$  is symmetric, then  $\hat{\beta}$  is also an unbiased estimator of  $\beta$ ,

$$E(\hat{\beta}) = \beta.$$

2. The sampling distribution of  $\hat{\beta}$ , when  $\Sigma_i$  is estimated from the data, is approximately multivariate normal with mean,  $\beta$ , and covariance

$$\text{Cov}(\hat{\beta}) = \left\{ \sum_{i=1}^N (X_i' \Sigma_i^{-1} X_i) \right\}^{-1}.$$

Furthermore these properties of  $\hat{\beta}$  hold in large samples even when the assumption that  $Y_i$  has a multivariate normal distribution is not valid, provided the data are complete. Thus an important property of the ML estimator of  $\beta$ , derived under the assumption of a multivariate normal distribution for  $Y_i$  given  $X_i$ , is that it provides a valid estimate of  $\beta$  even when the multivariate normal distribution assumption does not hold. Moreover this appealing property of the ML estimator of  $\beta$ , and of any GLS estimator of  $\beta$  (recall, the ML estimator of  $\beta$  is also the GLS estimator with the ML estimate of  $\Sigma_i(\theta)$  substituted), extends to the incomplete data setting when certain assumptions about missingness hold.

Thus, in terms of properties of the sampling distribution of  $\hat{\beta}$ , there is no penalty for actually having to estimate  $\Sigma_i$  from the longitudinal data at hand. However, as comforting as this result may appear to be, it must be kept in mind that this is a large sample (i.e., as  $N$  approaches infinity) property of  $\hat{\beta}$ . With sample sizes of the magnitude often encountered in many fields of application, the properties of the sampling distribution of  $\hat{\beta}$  can be expected to be adversely influenced by the estimation of a very large number of covariance parameters. This is an important issue that we will return to in Chapter 7.

## 4.3 MISSING DATA ISSUES

Although most longitudinal studies are designed to collect data on every individual in the sample at each time of follow-up, many studies have some missing observations. In longitudinal studies in the health sciences, missing data are the rule, not the exception. Missing data have three important implications for longitudinal analysis. First, when longitudinal data are missing, the data set is necessarily unbalanced over time since not all individuals have the same number of repeated measurements at a common set of occasions. As a result methods of analysis need to be able to handle the unbalanced data without having to discard data on individuals with any missing data. This feature of missingness will not be of any concern for the methods described in later chapters of the book. Second, when there are missing data, there will be a loss of information and a reduction in the precision with which changes in the mean response over time can be estimated. This reduction in precision is directly related to the amount of missing data and will also be influenced to a certain extent by how the analysis handles the missing data. For example, using only the complete cases (i.e., those individuals with no missing data) will usually be the least efficient method. Finally, when there are missing data, the validity of any method of analysis will require that certain assumptions about the reasons for any missingness, often referred to as the *missing data mechanism*, are tenable. Consequently, when data are missing we must carefully consider the reasons for missingness.

In this section we review two general types of missing data mechanisms. The two mechanisms differ in terms of assumptions concerning whether missingness is related to responses that have been observed. The distinctions between different types of missing data mechanisms and alternative methods for handling missingness in longitudinal studies will be discussed in greater detail in Chapters 17 and 18.

The missing data mechanism can be thought of as a model that describes the probability that a response is observed or missing at any occasion. We make an important distinction between missing data mechanisms that are referred to as *missing completely at random* (MCAR) and *missing at random* (MAR). The distinction between these two mechanisms determines the appropriateness of maximum likelihood estimation under the assumption of a multivariate normal distribution for the responses and GLS without requiring assumptions about the shape of the distribution.

Data are said to be missing completely at random (MCAR) when the probability that responses are missing is unrelated to either the specific values that, in principle, should have been obtained (the *missing responses*) or the set of observed responses. That is, longitudinal data are MCAR when missingness in  $Y_i$  is simply the result of a chance mechanism that does not depend on either observed or unobserved components of  $Y_i$ . The essential feature of MCAR is that the observed data can be thought of as a random sample of the complete data. As a result the moments (e.g., the means, variances, and covariances), and indeed, the distribution of the observed data do not differ from the corresponding moments or distribution of the complete data.

An MCAR mechanism has important consequences for the analysis of longitudinal data. In particular, any method of analysis that yields valid inferences in the absence of missing data will also yield valid inferences when missing data are MCAR and the analysis is based on all available data, or even when it is restricted to the “completers” (i.e., those with no missing data). Given that valid estimates of the means, variances, and covariances can be obtained, GLS provides valid estimates of  $\beta$  without requiring any distributional assumptions for  $Y_i$ . The GLS estimator of  $\beta$  is valid provided the model for the mean response has been correctly specified; it does not require any assumptions about the joint distribution of the longitudinal responses. The maximum likelihood (ML) estimator of  $\beta$ , under the assumption that the responses have a multivariate normal distribution, is also the GLS estimator (with the ML estimate of  $\sum_i(\theta)$ , e.g.,  $\hat{\Sigma}_i = \Sigma_i(\hat{\theta})$ , substituted). Thus in this setting the ML and GLS estimators have exactly the same properties regardless of the true distribution of  $Y_i$ .

In contrast to MCAR, data are said to be missing at random (MAR) when the probability that responses are missing depends on the set of observed responses, but is unrelated to the specific missing values that, in principle, should have been obtained. Put another way, if subjects are

stratified on the basis of similar values for the responses that have been observed, missingness is simply the result of a chance mechanism that does not depend on the values of the unobserved responses. However, because the missingness mechanism now depends on observed responses, the distribution of  $Y_i$  in each of the distinct strata defined by the patterns of missingness is not the same as the distribution of  $Y_i$  in the target population. This has important consequences for analysis. One is that an analysis restricted to the “completers” is not valid. Put another way, the “completers” are a biased sample from the target population. Furthermore the distribution of the observed components of  $Y_i$ , in each of the distinct strata defined by the patterns of missingness, does not coincide with the distribution of the same components of  $Y_i$  in the target population. Therefore the sample means, variances, and covariances based on either the “completers,” or the available data are biased estimates of the corresponding parameters in the target population. As a result GLS no longer provides valid estimates of  $\beta$  without making correct assumptions about the joint distribution of the longitudinal responses. On the other hand, ML estimation of  $\beta$  is valid when data are MAR provided that the multivariate normal distribution has been correctly specified. This requires correct specification of not only the model for the mean response but also the model for the covariance among the responses. In a sense, ML estimation allows the missing values to be validly “predicted” or “imputed” using the observed data and a correct model for the joint distribution of the responses.

To summarize, we have distinguished between two types of missing data mechanisms that are referred to as *missing completely at random* (MCAR) and *missing at random* (MAR). The MAR assumption is far less restrictive than MCAR. The distinction between these two mechanisms determines the appropriateness of maximum likelihood estimation under the assumption of a multivariate normal distribution for the responses and GLS without requiring assumptions about the shape of the distribution. The general properties of GLS described in the previous section require either that the data are complete or that any missing data are MCAR. If data are MAR, GLS based only on the means, variances, and covariances of the available data can yield biased estimates of  $\beta$ . In contrast, ML estimation yields valid estimates of  $\beta$  when data are MCAR or MAR, but for the latter mechanism, at the cost of requiring that the joint distribution of the responses is correctly specified. A more detailed discussion of missing data mechanisms, with concrete examples, and the implications of different types of missing data mechanisms for analysis, is presented in Chapter 17.

## 4.4 STATISTICAL INFERENCE

Next we consider how to make inferences about  $\beta$ . In particular, we consider the construction of confidence intervals and tests of hypotheses. To construct confidence intervals and tests of hypotheses about  $\beta$ , we can make direct use of the ML estimate  $\hat{\beta}$ , and its estimated covariance matrix

$$\widehat{\text{Cov}}(\hat{\beta}) = \left\{ \sum_{i=1}^N (X_i' \hat{\Sigma}_i^{-1} X_i) \right\}^{-1},$$

where  $\Sigma_i$  in (4.5) is replaced by  $\hat{\Sigma}_i$ , the ML estimate of  $\Sigma_i$ . For example, for any single component of  $\beta$ , say  $\beta_k$ , a natural method for constructing 95% confidence limits is by taking  $\hat{\beta}_k$  plus or minus 1.96 times the standard error of  $\hat{\beta}_k$ . Note that different confidence limits (e.g., 90%) can be obtained by choosing appropriate multiples of the standard error, based on the standard normal distribution. The standard error of  $\hat{\beta}_k$  is simply the square-root of the diagonal element of  $\widehat{\text{Cov}}(\hat{\beta})$  corresponding to  $\hat{\beta}_k$ ,

$$\sqrt{\widehat{\text{Var}}(\hat{\beta}_k)}.$$

Similarly a test of the null hypothesis,  $H_0: \beta_k = 0$  versus  $H_A: \beta_k \neq 0$ , can be based on the following Wald statistic:

$$Z = \frac{\hat{\beta}_k}{\sqrt{\widehat{\text{Var}}(\hat{\beta}_k)}},$$

where  $\widehat{\text{Var}}(\hat{\beta}_k)$  denotes the diagonal element of  $\widehat{\text{Cov}}(\hat{\beta})$  corresponding to  $\hat{\beta}_k$ . This test statistic can be compared with a standard normal distribution.

More generally, it may be of interest to construct confidence intervals and tests of hypotheses about certain linear combinations of the components of  $\beta$ . Let  $L$  denote a vector or matrix of *known* weights, and suppose that it is of interest to test  $H_0: L\beta = 0$ . The linear combination of the components of  $\beta$ ,  $L\beta$ , represents a contrast of scientific interest. For example, suppose that  $\beta = (\beta_1, \beta_2, \beta_3)'$  and let  $L = (0, 0, 1)$ , then  $H_0: L\beta = 0$  is equivalent to  $H_0: \beta_3 = 0$ . Alternatively, if  $L = (0, 1, -1)$ , then  $H_0: L\beta = 0$  is equivalent to  $H_0: \beta_2 - \beta_3 = 0$  or  $H_0: \beta_2 = \beta_3$ . A natural estimate of  $L\beta$  is given by  $L\hat{\beta}$ . Moreover it can be shown that the sampling distribution of  $L\hat{\beta}$  is multivariate normal with mean,  $L\beta$ , and with covariance matrix,  $L\widehat{\text{Cov}}(\hat{\beta})L'$ .

Note that in the two examples considered earlier,  $L$  is a single,  $1 \times 3$  row vector,  $L = (0, 0, 1)$  or  $L = (0, 1, -1)$ . If  $L$  is a single row vector, then  $L\widehat{\text{Cov}}(\hat{\beta})L'$  is a single value (or scalar) and its square-root provides an estimate of the standard error for  $L\hat{\beta}$ . Thus an approximate 95% confidence interval for  $L\beta$  is given by

$$L\hat{\beta} \pm 1.96 \sqrt{L\widehat{\text{Cov}}(\hat{\beta})L'}.$$

Similarly, in order to test  $H_0: L\beta = 0$  versus  $H_A: L\beta \neq 0$ , we can use the Wald statistic,

$$Z = \frac{L\hat{\beta}}{\sqrt{L\widehat{\text{Cov}}(\hat{\beta})L'}},$$

and compare this test statistic to a standard normal distribution. Recall that if  $Z$  is a standard normal random variable, then  $Z^2$  has a  $\chi^2$  distribution with 1 degree of freedom (df), denoted  $\chi^2_1$ . Thus an identical test of  $H_0: L\beta = 0$  versus  $H_A: L\beta \neq 0$ , uses the statistic

$$W^2 = (L\hat{\beta})\{L\widehat{\text{Cov}}(\hat{\beta})L'\}^{-1}(L\hat{\beta}),$$

and compares  $W^2$  to a  $\chi^2$  distribution with 1 degree of freedom. This latter observation helps to motivate how the Wald test readily generalizes when  $L$  has more than one row, thereby allowing simultaneous testing of a single multivariate hypothesis. For example, suppose that  $\beta = (\beta_1, \beta_2, \beta_3)'$  and it is of interest to test the equality of the three regression parameters. The null hypothesis can be expressed as  $H_0: \beta_1 = \beta_2 = \beta_3$ . Letting

$$L = \begin{pmatrix} 1 & -1 & 0 \\ 1 & 0 & -1 \end{pmatrix},$$

this null hypothesis can also be expressed as  $H_0: L\beta = 0$ , since if

$$\begin{aligned} L\beta &= \begin{pmatrix} 1 & -1 & 0 \\ 1 & 0 & -1 \end{pmatrix} \begin{pmatrix} \beta_1 \\ \beta_2 \\ \beta_3 \end{pmatrix} \\ &= \begin{pmatrix} \beta_1 - \beta_2 \\ \beta_1 - \beta_3 \end{pmatrix} = 0, \end{aligned}$$

then

$$\begin{pmatrix} \beta_1 \\ \beta_2 \\ \beta_3 \end{pmatrix} = \begin{pmatrix} \beta_1 \\ \beta_2 \\ \beta_3 \end{pmatrix},$$

or, equivalently,  $\beta_1 = \beta_2 = \beta_3$ . In general, suppose that  $L$  has  $r$  rows (e.g., representing  $r$  contrasts of scientific interest), then a simultaneous test of  $H_0: L\beta = 0$  versus  $H_A: L\beta \neq 0$  is given by

$$W^2 = (L\hat{\beta})' \{L\widehat{\text{Cov}}(\hat{\beta})L'\}^{-1} (L\hat{\beta}),$$

which has a  $\chi^2$  distribution with  $r$  df. The latter test is often referred to as a multivariate Wald test.

One alternative to the Wald test is the *likelihood ratio test*. The likelihood ratio test of  $H_0: L\beta = 0$  versus  $H_A: L\beta \neq 0$  is obtained by comparing the maximized log-likelihoods for two models, one model that incorporates the constraint that  $L\beta = 0$  (e.g.,  $\beta_3 = 0$  or  $\beta_2 = \beta_3$ ), the other model unconstrained (i.e., without the constraint,  $L\beta = 0$ ). The latter is referred to as the “full” model and the former is referred to as the “reduced” model. Note that these two models are *nested*, in the sense that the “reduced” model is a special case of the “full” model. That is, when the reduced model is *nested* within the full model, it is a particular version of the full model, so that when the reduced model holds, the full model must necessarily hold.

The likelihood ratio test for two nested models can be constructed by comparing their respective maximized log-likelihoods, say  $\hat{l}_{\text{full}}$  and  $\hat{l}_{\text{red}}$ , for the full and reduced models, respectively. The former is at least as large as the latter. The larger the difference between  $\hat{l}_{\text{full}}$  and  $\hat{l}_{\text{red}}$ , the stronger the evidence is that the reduced model is inadequate. A formal statistical test is obtained by taking twice the difference in the respective maximized log-likelihoods,

$$G^2 = 2(\hat{l}_{\text{full}} - \hat{l}_{\text{red}}),$$

and comparing the statistic to a chi-squared distribution with degrees of freedom equal to the difference between the number of parameters in the full and reduced models. This test is called the *likelihood ratio test*. We caution that the likelihood ratio test can only be used when the number of observations for the “full” and “reduced” models is the same. For example, when there are some missing data on a covariate that belongs to the “full” model only, the “full” and “reduced” models are no longer nested. Instead, to ensure the models are nested, the likelihood ratio test must be based on a comparison of the “full” and “reduced” models for the same subset of observations with no missing values on that covariate.

This use of the likelihood can also provide confidence limits for  $\beta$  or  $L\beta$ . Rather than calculating confidence limits for  $\beta$  (or  $L\beta$ ) as the maximum likelihood estimate,  $\hat{\beta}$ , plus or minus an appropriate multiple of the standard errors, likelihood-based confidence intervals can be constructed. The basic idea behind likelihood-based confidence intervals is to consider all values of  $\beta$  (or  $L\beta$ ) that are consistent with the data at hand. More formally, for a single component of  $\beta$ , say  $\beta_k$ , we can define a *profile log-likelihood*,  $l_p(\beta_k)$ , obtained by maximizing the log-likelihood over the remaining parameters while holding  $\beta_k$  at some fixed value. A likelihood-based confidence interval for  $\beta_k$  is obtained by considering values of  $\beta_k$  that are reasonably consistent with the data. Specifically, an approximate 95% likelihood-based confidence interval is given by the set of all values of  $\beta_k$  satisfying

$$2 \times \{l_p(\hat{\beta}_k) - l_p(\beta_k)\} \leq 3.84,$$

where the critical value on the right-hand side of the equation is obtained from a chi-squared distribution with 1 degree of freedom. More generally, confidence intervals for  $L\beta$  can be obtained by inverting the corresponding test of  $H_0: L\beta = 0$  in a similar way.

Although the construction of likelihood ratio tests and likelihood-based confidence intervals is more involved (e.g., requiring an additional fit of the model under the null hypothesis) than the corresponding Wald-based tests and confidence intervals, the likelihood-based tests and confidence intervals often have superior properties. This is especially the case when the response variable is discrete. For example, in logistic regression with binary data, likelihood ratio tests have better properties than the corresponding Wald tests. Thus, when in doubt, we recommend the use of likelihood-based tests and confidence intervals. However, for ease of presentation, many of the results presented in later chapters rely on Wald-based tests and confidences intervals; likelihood-based tests and confidence intervals are presented only in cases where the discrepancies might change the substantive conclusions of the analysis.

Finally, we note that likelihood ratio tests can also be used for hypotheses about the covariance parameters. However, there are some potential problems with the standard use of the likelihood ratio test for comparing nested models for the covariance; we will return to this topic in Chapter 7. In general, we do not recommend testing hypotheses about the covariance parameters using Wald tests (i.e., based on the ratio of the parameter estimate to its standard error). In particular, the sampling distribution of the Wald test statistic for a variance parameter does not have an approximate normal distribution when the sample size is relatively small and the population variance is close to zero. Because the variance has a lower bound of zero, very large samples are required to justify the normal approximation for the sampling distribution of the Wald test statistic when the variance is close to zero.

# Comment on Denominator Degrees of Freedom

So far in our discussion of confidence intervals and tests of hypotheses about  $\beta$  we have relied on the large sample properties of the sampling distribution of the ML estimate of  $\beta$ . That is, we have used the standard normal and chi-squared distributions instead of  $t$  and  $F$  distributions. It can be argued that the use of the standard normal and chi-squared distributions is more “liberal” (or “anti-conservative”) than the corresponding  $t$  and  $F$  distributions because there is an implicit assumption of infinite denominator degrees of freedom. By “liberal,” we mean that nominal  $p$ -values may be too small and confidence intervals may be too narrow.

With large denominator degrees of freedom (e.g., due to a large sample size) estimation of  $\theta$  or  $\Sigma_i$  does not introduce any additional uncertainty. However, with small sample sizes there is some uncertainty attached to the estimation of  $\theta$  that needs to be acknowledged in our inferences about  $\beta$ . Ordinarily this additional source of uncertainty is recognized by use of the  $t$  and  $F$  distributions instead of the standard normal and chi-squared distributions.

A practical difficulty with the use of the  $t$  and  $F$  distributions in this setting is that the denominator degrees of freedom associated with tests and confidence intervals for components of  $\beta$  is not easy to determine except in certain special cases where the data are balanced and the model for the mean has a relatively simple form. To circumvent this difficulty, various approximations for the denominator degrees of freedom have been proposed. One well-known method is the Satterthwaite approximation, a somewhat tedious and computationally demanding procedure. If Satterthwaite’s (1946) method is used to obtain approximate denominator degrees of freedom, say  $\hat{v}$ , then an approximate 95% confidence interval for  $L\beta$  is given by

$$L\hat{\beta} \pm t_{\hat{v}, 0.025} \sqrt{L \widehat{\text{Cov}}(\hat{\beta}) L'}$$

where  $t_{\hat{v}, 0.025}$  is the upper 2.5% cutoff from a  $t$  distribution with  $\hat{v}$  degrees of freedom (i.e., for the  $t$  distribution with  $\hat{v}$  degrees of freedom, 95% of the area lies between  $-t_{\hat{v}, 0.025}$  and  $t_{\hat{v}, 0.025}$ ). Similarly, to test  $H_0: L\beta = 0$  versus  $H_A: L\beta \neq 0$ , we can use the Wald statistic

$$\frac{L\hat{\beta}}{\sqrt{L \widehat{\text{Cov}}(\hat{\beta}) L'}}$$

and compare this test statistic to a  $t$  distribution with  $\hat{v}$  degrees of freedom. The Satterthwaite approximation can also be applied to multivariate Wald statistics, with the chi-squared distribution replaced by the  $F$  distribution (when the multivariate Wald statistic has been divided by the number of rows of the matrix  $L$  or the numerator degrees of freedom).

Recently Kenward and Roger (1997) proposed an alternative approximation that adjusts the test statistics and provides approximate denominator degrees of freedom. Although the Satterthwaite approximation and the approximation proposed by Kenward and Roger (1997) are implemented as options in some statistical software packages (e.g., PROC MIXED in SAS), it must be emphasized that the small sample properties of these approximations in regression models for longitudinal data have not been extensively studied.

In summary, the use of the standard normal and chi-squared distributions is valid when  $\Sigma_i$  (or  $\theta$ ) is known, or when  $\Sigma_i$  has been estimated with a large number of degrees of freedom. Recall that there is not much practical difference between the use of the standard normal and  $t$  distributions once the degrees of freedom of the latter exceed 100. With small sample sizes, there is some uncertainty in the estimation of  $\theta$  that should be accounted for and the use of the  $t$  and  $F$  distributions, with degrees of freedom approximated by the methods of Satterthwaite (1946) or Kenward and Roger (1997), is preferred. Fortunately, in many applications in the health sciences the numbers of subjects is reasonably large relative to the number of measurement occasions. As a result the unknown denominator degrees of freedom, especially for components of  $\beta$  that represent time trends and their interactions with covariates (e.g., group  $\times$  time interactions), will be sufficiently large that the standard normal and chi-squared distributions are reasonable approximations to the corresponding  $t$  and  $F$  distributions. For the remainder of the book, we construct confidence intervals and tests of

hypotheses about  $\beta$  using the standard normal and chi-squared distributions; we use approximations for the denominator degrees of freedom only in cases where it might change the substantive conclusions of the analysis.

## 4.5 RESTRICTED MAXIMUM LIKELIHOOD (REML) ESTIMATION

We conclude this chapter with a discussion of a variant on ML estimation, known as *restricted maximum likelihood* (REML) estimation. Recall that the ML estimates of  $\beta$  and  $\theta$  (or  $\Sigma_i$ ) were obtained by maximizing the following log-likelihood function:

$$l = -\frac{K}{2} \log(2\pi) - \frac{1}{2} \sum_{i=1}^N \log |\Sigma_i| - \frac{1}{2} \left\{ \sum_{i=1}^N (\mathbf{y}_i - \mathbf{X}_i \boldsymbol{\beta})' \Sigma_i^{-1} (\mathbf{y}_i - \mathbf{X}_i \boldsymbol{\beta}) \right\}.$$

Although the ML estimates of  $\beta$  and  $\Sigma_i$  ( $\theta$ ) have desirable large sample (or asymptotic) properties, the ML estimate of  $\Sigma_i$  has a well-known bias in finite samples. For example, the diagonal elements of  $\Sigma_i$  are underestimated.

To illustrate the problem, consider the case where data arise from a series of cross-sectional studies that are repeated at  $n$  different occasions. Here we can assume that the observations are independent of one another, and for ease of exposition we also assume that the variance is constant, say  $\sigma^2$ . As noted earlier, the ML estimates of  $\beta$  and  $\sigma^2$  are obtained by maximizing

$$-\frac{K}{2} \log(2\pi\sigma^2) - \frac{1}{2} \sum_{i=1}^N \sum_{j=1}^n (y_{ij} - X'_{ij} \boldsymbol{\beta})^2 / \sigma^2.$$

The ML estimator of  $\beta$  is

$$\hat{\boldsymbol{\beta}} = \left\{ \sum_{i=1}^N \sum_{j=1}^n (X_{ij} X'_{ij}) \right\}^{-1} \sum_{i=1}^N \sum_{j=1}^n (X_{ij} y_{ij}),$$

while the ML estimator of  $\sigma^2$  is

$$\hat{\sigma}^2 = \sum_{i=1}^N \sum_{j=1}^n (y_{ij} - X'_{ij} \hat{\boldsymbol{\beta}})^2 / K,$$

where  $K = n \times N$ . Furthermore, it can be shown that

$$E(\hat{\sigma}^2) = \left( \frac{K-p}{K} \right) \sigma^2,$$

where  $p$  is the dimension of  $\beta$ . As a result, the ML estimate of  $\sigma^2$  is biased in finite samples and underestimates  $\sigma^2$ . An unbiased estimator is obtained by using  $K-p$  (or the residual degrees of freedom) as the denominator instead of  $K$ ,

$$\tilde{\sigma}^2 = \sum_{i=1}^N \sum_{j=1}^n (y_{ij} - X'_{ij} \hat{\boldsymbol{\beta}})^2 / (K-p).$$

This estimator for  $\sigma^2$  is also known as the REML estimator. Note that the bias of the ML estimate of  $\sigma^2$  is a decreasing function of the total number of observations,  $K$ .

In effect, the bias arises because the ML estimate has not taken into account the fact that  $\beta$  is also estimated from the data. In the estimator of  $\sigma^2$  we have replaced  $\beta$  by  $\hat{\boldsymbol{\beta}}$  but have failed to acknowledge in some sense that  $\beta$  was estimated from the data. If there are problems of bias with the ML estimate of  $\sigma^2$  with independent observations, then it should not come as a great surprise that similar problems arise in the estimation of  $\Sigma_i$  (or  $\theta$ ) with correlated data.

The theory of restricted (or residual) maximum likelihood (REML) estimation was developed to address this problem. The main idea behind REML estimation is to separate that part of the data used for estimation of  $\Sigma_i$  from that used for estimation of  $\beta$ . Estimation of  $\Sigma_i$  is then based only on the relevant part of the data. Thus the fundamental idea in REML estimation of  $\Sigma_i$  is to eliminate  $\beta$  from the likelihood so that it is defined only in terms of  $\Sigma_i$ . This can be achieved in a number of ways. One possible way to obtain the restricted likelihood is to transform the data to a set of linear combinations of observations that have a distribution that does not depend on  $\beta$ . For example, the residuals after estimating  $\beta$  by ordinary least squares (OLS) can be used as the data for estimating  $\Sigma_i$  (or  $\theta$ ). The likelihood for these residuals depends only on  $\theta$ , and not on  $\beta$ . Thus, rather than maximizing the log-likelihood

$$(4.7) \quad -\frac{1}{2} \sum_{i=1}^N \log |\Sigma_i| - \frac{1}{2} \sum_{i=1}^N (y_i - X_i \hat{\beta})' \Sigma_i^{-1} (y_i - X_i \hat{\beta}),$$

REML maximizes the following slightly modified log-likelihood (formed from the residuals)

$$(4.8) \quad \begin{aligned} -\frac{1}{2} \sum_{i=1}^N \log |\Sigma_i| & - \frac{1}{2} \sum_{i=1}^N (y_i - X_i \hat{\beta})' \Sigma_i^{-1} (y_i - X_i \hat{\beta}) \\ & - \frac{1}{2} \log \left| \sum_{i=1}^N X_i' \Sigma_i^{-1} X_i \right|. \end{aligned}$$

When the residual likelihood given by (4.8) is maximized, we obtain an estimate of  $\theta$  (or  $\Sigma_i(\theta)$ ) that has made a correction for the fact that  $\beta$  has also been estimated. Of note, the additional term in the REML log-likelihood involves a determinant term,

$$\begin{aligned} -\frac{1}{2} \log \left| \sum_{i=1}^N X_i' \Sigma_i^{-1} X_i \right| & = \frac{1}{2} \log \left| \left( \sum_{i=1}^N X_i' \Sigma_i^{-1} X_i \right)^{-1} \right| \\ & = \log \left| \text{Cov}(\hat{\beta}) \right|^{\frac{1}{2}}, \end{aligned}$$

that can be expressed as the covariance of  $\hat{\beta}$ . As a result the REML likelihood multiplies the usual ML likelihood by a factor that is the square-root of the *generalized variance* of  $\hat{\beta}$ , a single number summary of the variation in the estimate of  $\beta$ . This makes a correction or adjustment that is analogous to the correction to the denominator in  $\hat{\sigma}^2$ .

We recommend the use of the REML estimator for  $\Sigma_i$ . In general, the REML estimator will be less seriously biased than the ML estimator for  $\Sigma_i$ . It should be noted that the difference between ML and REML estimation becomes less important when the sample size,  $N$ , is substantially larger than  $p$ , the dimension of  $\beta$ . Finally, when REML estimation is used to estimate  $\Sigma_i$ ,  $\beta$  is estimated by the usual generalized least squares (GLS) estimator

$$\hat{\beta} = \left\{ \sum_{i=1}^N (X_i' \hat{\Sigma}_i^{-1} X_i) \right\}^{-1} \sum_{i=1}^N (X_i' \hat{\Sigma}_i^{-1} y_i),$$

where  $\hat{\Sigma}_i = \Sigma_i(\hat{\theta})$  is the REML estimate of  $\Sigma_i$ .

On a final note, while the REML log-likelihood can be used to compare nested models for the covariance (e.g., in terms of likelihood ratio tests comparing nested models for the covariance), it should not be used to compare nested regression models for the mean. The extra determinant term in the REML log-likelihood depends on the regression model specification. As a result the REML likelihoods for two nested models for the mean response are based on quite different transformations of the data (to obtain linear combinations of  $y_i$  whose distributions do not depend on  $\beta$ ). In short, the REML likelihoods for two nested models for the mean are based on two entirely different sets of transformed responses, making comparisons between the models meaningless. Instead, the standard ML log-likelihood should be used for constructing likelihood ratio tests that compare nested regression models for the mean.

In conclusion, we recommend the use of REML for estimation of  $\Sigma_i$  (with  $\beta$  estimated using the GLS estimator that substitutes the REML estimate,  $\hat{\Sigma}_i$ , for  $\Sigma_i$ ). The REML log-likelihood should also be used in comparing nested models for the covariance. However, the construction of likelihood ratio tests comparing nested models for the mean should always be based on the ML, not the REML, log-likelihood.

## **4.6 FURTHER READING**

Many textbooks on statistical theory and methods include a discussion of the methods of least squares and maximum likelihood estimation. Weisberg (1985) provides a useful introduction to the method of least squares in the context of regression; Chapter 4 of Cox and Wermuth (1996) presents a concise but remarkably lucid description of least squares, generalized least squares, and maximum likelihood estimation.

## Bibliographic Notes

A discussion of the properties of generalized least squares (GLS) estimators can be found in, for example, Amemiya (1985) and Newey and McFadden (1994). Kakwani (1967) and Kackar and Harville (1981) discuss the unbiasedness properties of GLS estimators when the assumption of normally distributed errors is replaced by the weaker assumption that the distribution of the errors is symmetric.

The use of REML, as an alternative to maximum likelihood, for covariance parameter estimation was originally proposed by Patterson and Thompson (1971). Special cases of REML estimation had previously been considered by Anderson and Bancroft (1952), Russell and Bradley (1958), and Thompson (1962) in the context of balanced ANOVA models. Harville (1974) presented a Bayesian interpretation of REML.

# *Chapter 5*

## *Modeling the Mean: Analyzing Response Profiles*

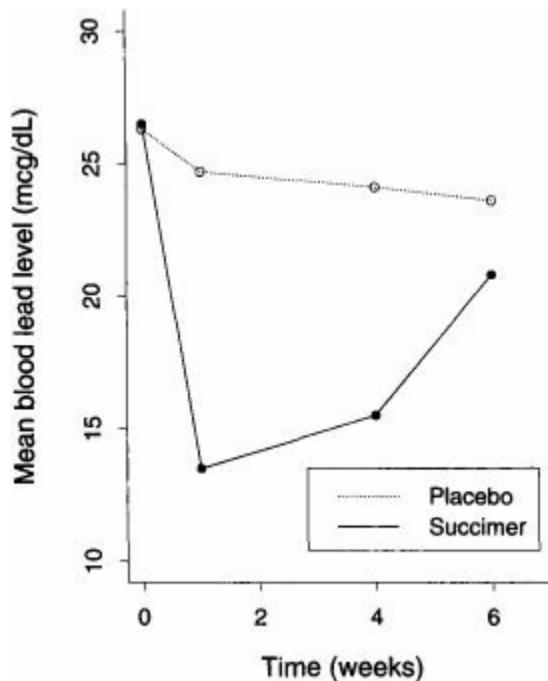
### **5.1 INTRODUCTION**

In this chapter we present a method for analyzing longitudinal data that imposes minimal structure or restrictions on the mean response over time and on the covariance among the repeated measures. The method focuses on analyzing response profiles and can be applied to longitudinal data when the design is balanced, with the timing of the repeated measures common to all individuals in the study. Although we focus on study designs where all subjects are measured at the same set of  $n$  occasions (i.e., *balanced* longitudinal designs), as we show, the analysis of response profiles can also handle incompleteness due to missing data (i.e., incomplete longitudinal studies with balanced designs).

Methods for analyzing response profiles are appealing when there is a single categorical covariate (perhaps denoting different treatment or exposure groups) and when no specific *a priori* pattern for the differences in the response profiles between groups can be specified. When repeated measures are obtained at the same sequence of occasions, the data can be summarized by the estimated mean response at each occasion, stratified by levels of the group factor. At any given level of the group factor, the sequence of means over time is referred to as the mean *response profile*.

For example, consider the blood lead level data from the TLC trial. The mean response profiles for the two groups randomized to succimer and placebo are presented in [Figure 5.1](#). This plot is produced by simply calculating the arithmetic average of the responses at each occasion, within each treatment group, and joining adjacent means with a series of line segments. In settings where data on some subjects are missing, such a plot can still be made but it is obtained from the estimated means for each occasion, stratified by group. We will show how to estimate these mean response profile curves in Section 5.3.

[\*\*Fig. 5.1\*\*](#) Mean blood lead levels at baseline (week 0), week 1, week 4, and week 6 in the succimer and placebo groups.



The main goal in the analysis of response profiles is to characterize the patterns of *change* in the mean response over time in the groups and to determine whether the shapes of the mean response profiles differ among the groups. For example, in the TLC trial, the major question of scientific

interest is concerned with whether *changes* in the mean blood lead levels are the same for the succimer and placebo groups. In Sections 5.2 and 5.3 we show how questions about whether the patterns of change are the same in all groups translate into hypotheses about the interaction between the group factor and time.

Methods for analyzing response profiles can be extended in a straightforward way to handle the case where there is more than a single group factor and when there are baseline covariates that need to be adjusted for. However, for ease of exposition, we focus on the case where there is only a single group factor. For example, in an observational study the groups might be defined by characteristics of the study subjects, such as age, gender, or exposure level. Alternatively, groups might be defined by random assignment to different treatments or interventions. The distinction between observational studies and randomized trials is important and, as we will see later, has ramifications for the analysis of response profiles.

A characteristic feature of longitudinal studies is the presence of a baseline measurement. In the TLC trial, the objective is to compare the patterns of change in blood lead levels from baseline over time across the treatment groups. The baseline measurement is an outcome like those measured subsequently, but is unique in that, being pre-randomization, it can be assumed not to depend on treatment group. Indeed, this is apparent in the plot of the mean response profiles in [Figure 5.1](#). This is a common feature of longitudinal studies which involve randomization after baseline. The baseline response may play a special role in other settings as well. For example, sometimes the baseline response is range restricted, as when only subjects with values greater than or less than a threshold are included in the study. With observational studies of growth or decline, groups may be known to differ at baseline, or comparison groups may be selected by matching so that baseline means are comparable.

Thus the question naturally arises as to how to handle the baseline measurement in the assessment of change. This is important, since it will affect how we construct hypothesis tests, and how they should be interpreted. In addition, how we handle the baseline response in the analysis will have an impact on efficiency and the power of tests of hypotheses. In Section 5.6, we describe two ways of adjusting for the baseline value in a simple setting and discuss their relative merits under different longitudinal study designs. In Section 5.7, we compare and contrast a number of alternative strategies for handling the baseline response in more general settings and make recommendations about the preferred strategies in different situations. Many of our readers may find the level of detail in Section 5.7 somewhat daunting. We note that Section 5.7 can be omitted at first reading without loss of continuity. However, we encourage all of our readers to eventually tackle the material in Section 5.7 since appropriate adjustment for baseline is an important aspect of the analysis of longitudinal change.

## 5.2 HYPOTHESES CONCERNING RESPONSE PROFILES

In our discussion of the analysis of response profiles, we focus initially on the two-group design, but generalizations to more than two groups are straightforward. Given a sequence of  $n$  repeated measures on a number of distinct groups of individuals, three main questions concerning the response profiles can be posed:

1. Are the mean response profiles similar in the groups, in the sense that the mean response profiles are parallel?

This is a question that concerns the *group × time interaction effect*. A graphical representation of the null hypothesis of parallel mean response profiles is displayed in [Figure 5.2\(a\)](#).

2. Assuming that the population mean response profiles are parallel, are the means constant over time, in the sense that the mean response profiles are flat?

This is a question that concerns the *time effect*. A graphical representation of the null hypothesis that the mean response profiles are flat is displayed in [Figure 5.2\(b\)](#).

3. Assuming that the population mean response profiles are parallel, are they also at the same level in the sense that the mean response profiles for the groups coincide?

This is a question that concerns the *group effect*. A graphical representation of the null hypothesis that the mean response profiles are at the same level is displayed in [Figure 5.2\(c\)](#).

In longitudinal studies, the first question is of main scientific interest. The hypothesis of parallel response profiles corresponds to the hypothesis that the patterns of change in the mean response over time are the same across groups. This comparison of change in the response over time is the raison d'être of a longitudinal study. In contrast, as we will see later, the second and third questions may not have any scientific relevance. That is, even when the response profiles can be assumed to be parallel, any interest in the second and third questions is secondary and depends on the longitudinal study design.

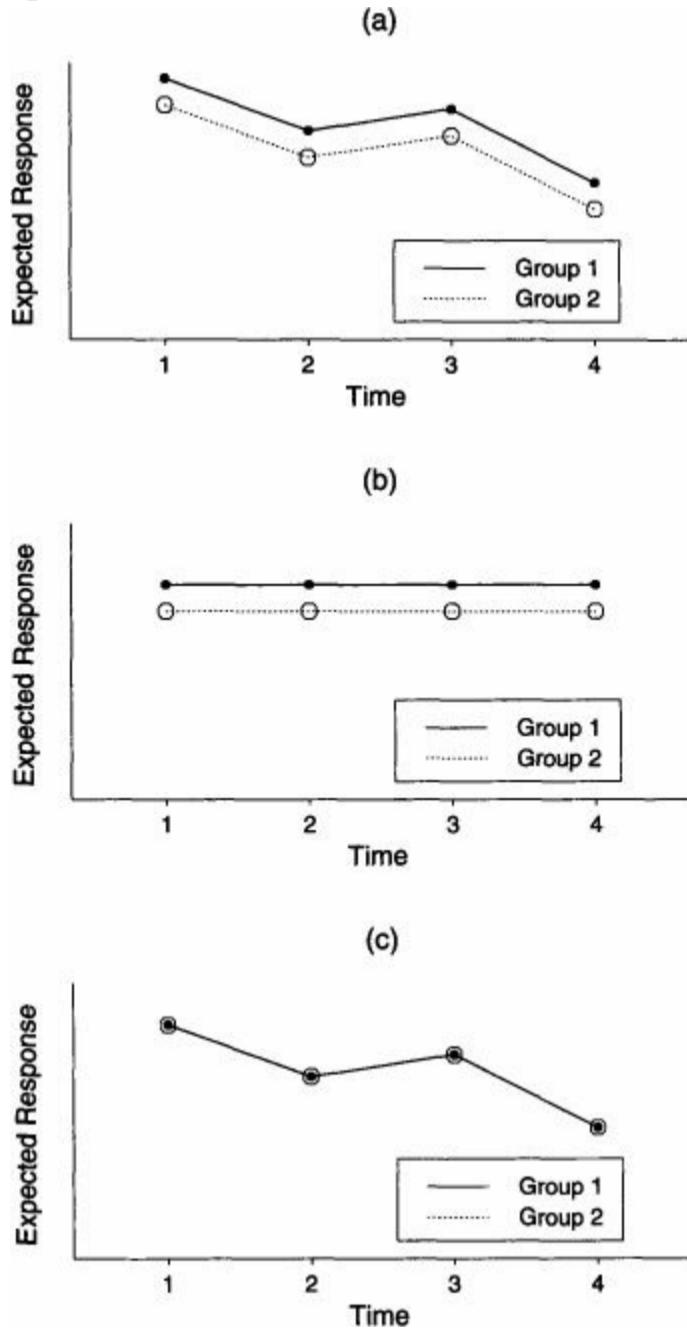
Note that the second and third questions have implicitly made an assumption about the answer to the first. There is a very good reason for doing so. Except in very rare circumstances, it is not meaningful to ask the second and third questions if the mean response profiles are not parallel. Indeed, this is consistent with the general principle that *main effects* (e.g., group or time effects) are ordinarily not of interest when there is an interaction among the factors. That is, when there is a  $\text{group} \times \text{time}$  interaction, the mean response profiles in the groups are different (non-parallel profiles), consequently their shape can be described only with reference to a specific group, and their level can be described only with reference to a specific time.

The appropriate scientific hypotheses in any particular study must be derived from the relevant scientific issues in that investigation. Here it becomes important to distinguish between longitudinal data arising from a randomized trial and from an observational study. In the former case, when study participants have been randomized to treatment groups and the baseline value of the response has been obtained prior to any study interventions, the mean response at occasion 1 is independent of treatment assignment. That is, by design, the group means are equal at baseline (occasion 1). In contrast, in an observational study, there is no a priori reason to assume that the groups have the same mean response at baseline unless the groups were selected by matching on baseline response.

Consider a randomized longitudinal clinical trial comparing treatments where the measurement at the first occasion is a baseline response, obtained prior to any study interventions. For example, in the TLC trial, the blood lead levels at baseline were obtained prior to receiving placebo or succimer. In that case the only question of scientific interest is the first because it addresses whether the patterns of change in the mean response over time are the same in all groups. For example, in the TLC trial, the test of the  $\text{group} \times \text{time}$  interaction assesses whether changes in the mean blood levels are the same for the succimer and placebo groups. In a randomized trial, the second question is usually of less importance because it does not involve a direct comparison of groups. The second

question concerns the time effect, where the focus is on the comparison of the mean response at each occasion averaged over the groups. Hypotheses concerning the main effect of time translate into questions concerning whether the overall (i.e., averaged over groups) mean response has changed from baseline. Finally, in a randomized trial where the baseline response has been obtained before any study interventions, there is no interest in the third question. The third question concerns the group effect. However, in this setting the absence of a group  $\times$  time interaction implies that there is no group effect. That is, if the groups have the same pattern of change over time and, by design, do not differ at baseline, their mean response profiles must necessarily coincide. As a result the test of group effect is subsumed within the test of group  $\times$  time interaction.

**Fig. 5.2** Graphical representation of the null hypotheses of (a) no group  $\times$  time interaction effect, (b) no time effect, and (c) no group effect.



In an observational study where the group factor might represent different exposures or inherent characteristics of the individuals, the first question is usually of primary interest. It addresses the fundamental question of whether patterns of change over time in the mean response vary by group. In contrast to a randomized trial, however, the second and third questions may also be of substantive interest. For example, in a longitudinal study of growth or aging, there may be interest in the pattern of change in the mean response over time, even when the pattern of change is the same in all groups. This concerns the main effect of time and is addressed by the second question. Ordinarily, when there is interest in the time effect, the preferred way to describe the trend in the mean response is with a relatively simple parametric curve; methods for fitting parametric or semiparametric curves to the mean response are described in Chapter 6. Finally, in an observational study, there may be interest in group comparisons of the mean response averaged over time. This concerns the group effect and is addressed by the third question. However, in the absence of any group  $\times$  time interaction, it must be recognized that the test for the main effect of group represents a comparison of the groups in terms of

their baseline (occasion 1) response. That is, if the groups have the same pattern of change over time, any group differences in the overall (i.e., averaged over occasions) mean response must reflect existing baseline differences among the groups.

To highlight the main features of the analysis of response profiles, consider the following example from a two-group study comparing a novel *treatment* and a *control*. We assume that the two groups have repeated measurements at the same set of  $n$  occasions. The analysis of response profiles is based on comparing the mean response profiles in the two groups. In a somewhat relaxed notation, let  $\mu(T) = \{\mu_1(T), \dots, \mu_n(T)\}'$  denote the mean response profile for the treatment group and  $\mu(C) = \{\mu_1(C), \dots, \mu_n(C)\}'$  denote the mean response profile for the control group. The population means in the two groups at each occasion are given in [Table 5.1](#).

**Table 5.1** Mean response profile over time in the treatment and control groups.

Group	Measurement Occasion			
	1	2	...	$n$
Treatment	$\mu_1(T)$	$\mu_2(T)$	...	$\mu_n(T)$
Control	$\mu_1(C)$	$\mu_2(C)$	...	$\mu_n(C)$
Difference	$\Delta_1$	$\Delta_2$	...	$\Delta_n$

Note:  $\Delta_j = \mu_j(T) - \mu_j(C)$ .

In this hypothetical study we are primarily interested in testing scientific hypotheses that compare the novel treatment and the control in terms of *changes* in the mean response over time. This can be determined by considering the null hypothesis of no group  $\times$  time interaction; that is, the null hypothesis that the mean response profiles are parallel. If the mean response profiles are parallel, the difference in the means between the two groups is constant over time. As a result in terms of  $\Delta_j$ , the null hypothesis is

$$H_{01}: \Delta_1 = \Delta_2 = \dots = \Delta_n,$$

where  $\Delta_j = \mu_j(T) - \mu_j(C)$ . If the null hypothesis is rejected, the two groups have non-parallel mean response profiles and the patterns of change over time differ in the two groups. Note that the number of constraints on the mean responses under this null hypothesis is  $n - 1$ . As a result the test of this null hypothesis has  $n - 1$  degrees of freedom. In Section 5.3 we describe how the constraints can be expressed in terms of specific contrasts of the means.

The previous illustration focused on the special case of  $G = 2$  groups; however, the main ideas can be generalized in a straightforward way when there are more than two groups. When there are  $G$  groups with repeated measurements at the same set of  $n$  occasions, we let  $\mu(g) = \{\mu_1(g), \dots, \mu_n(g)\}'$  denote the mean response profile for the  $g^{th}$  group ( $g = 1, \dots, G$ ). The population means in the  $G$  groups at each occasion are given in [Table 5.2](#). With  $G > 2$ , we can compare groups in a number of different ways. However, with  $G$  groups, there are only  $G - 1$  non-redundant comparisons. We define  $\Delta_j(g) = \mu_j(g) - \mu_j(G)$ , (for  $j = 1, \dots, n$ ;  $g = 1, \dots, G - 1$ ). That is,  $\Delta_j(g)$  is a contrast or comparison of the mean response at the  $j^{th}$  occasion for the  $g^{th}$  group (for  $g = 1, \dots, G - 1$ ) with the mean response at the  $j^{th}$  occasion in group  $G$ . Then the null hypothesis that the mean response profiles are parallel is

**Table 5.2** Mean response profile over time in  $G$  groups.

Group	Measurement Occasion			
	1	2	...	n
1	$\mu_1(1)$	$\mu_2(1)$	...	$\mu_n(1)$
2	$\mu_1(2)$	$\mu_2(2)$	...	$\mu_n(2)$
$\vdots$	$\vdots$	$\vdots$		$\vdots$
$g$	$\mu_1(g)$	$\mu_2(g)$	...	$\mu_n(g)$
$\vdots$	$\vdots$	$\vdots$		$\vdots$
$G$	$\mu_1(G)$	$\mu_2(G)$	...	$\mu_n(G)$

$$H_{01}: \Delta_1(g) = \Delta_2(g) = \dots = \Delta_n(g); \text{ for } g = 1, \dots, G - 1.$$

With  $G \geq 2$ , the test of the null hypothesis of no group  $\times$  time interaction effect has  $(G - 1) \times (n - 1)$  degrees of freedom.

So far our discussion of the analysis of response profiles has focused on an omnibus test of the group  $\times$  time interaction. However, unless the test of the group  $\times$  time interaction has only a single degree of freedom, this test does not help in discerning in what manner the patterns of change over time differ across groups. For the latter we must consider estimates (and their standard errors) of relevant contrasts of the means. In the next section, we describe a general linear model formulation of the analysis of response profiles. As we will show, the analysis of response profiles can be formulated in such a way that certain regression parameters have interpretations that bear directly on the questions of main scientific interest.

## 5.3 GENERAL LINEAR MODEL FORMULATION

Before we illustrate the main ideas with a numerical example, we consider how the analysis of response profiles can be implemented in the general linear model

$$E(Y_i|X_i) = \mu_i = X_i\beta,$$

for appropriate choices of  $X_i$ . We also describe how the main hypothesis of no group  $\times$  time interaction effect can be expressed in terms of  $\beta$ . Let  $n$  be the number of repeated measures and  $N$  the number of subjects. To express the model for the longitudinal design with  $G$  groups and  $n$  occasions of measurement, we require  $G \times n$  parameters for the  $G$  mean response profiles.

For example, with two groups measured at three occasions, there are  $2 \times 3 = 6$  mean parameters (see [Table 5.2](#)). For the first group, let the design matrix  $X_i$  be the following  $3 \times 6$  matrix:

$$X_i = \begin{pmatrix} 1 & 0 & 0 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 & 0 & 0 \end{pmatrix},$$

while for the second group, let the design matrix be

$$X_i = \begin{pmatrix} 0 & 0 & 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 0 & 0 & 1 \end{pmatrix}.$$

Then, in terms of the model

$$E(Y_i|X_i) = \mu_i = X_i\beta,$$

where  $\beta = (\beta_1, \dots, \beta_6)'$  is a  $6 \times 1$  vector of regression coefficients,

$$\mu(1) = \begin{pmatrix} \mu_1(1) \\ \mu_2(1) \\ \mu_3(1) \end{pmatrix} = \begin{pmatrix} \beta_1 \\ \beta_2 \\ \beta_3 \end{pmatrix};$$

similarly

$$\mu(2) = \begin{pmatrix} \mu_1(2) \\ \mu_2(2) \\ \mu_3(2) \end{pmatrix} = \begin{pmatrix} \beta_4 \\ \beta_5 \\ \beta_6 \end{pmatrix}.$$

As a result hypotheses about the mean response profiles in the two groups that were previously expressed in terms of  $\mu(1) = \{\mu_1(1), \mu_2(1), \mu_3(1)\}'$  and  $\mu(2) = \{\mu_1(2), \mu_2(2), \mu_3(2)\}'$  can easily be re-expressed in terms of hypotheses about the components of  $\beta$ . Specifically, the hypothesis of no group  $\times$  time interaction effect can be expressed as

$$H_{01}: (\beta_1 - \beta_4) = (\beta_2 - \beta_5) = (\beta_3 - \beta_6).$$

In this parameterization, hypotheses about the group  $\times$  time interaction cannot be expressed in terms of certain components of  $\beta$  being zero; instead, these hypotheses can be expressed in terms of  $L\beta = 0$ , for particular choices of vectors or matrices  $L$ . For example, the null hypothesis of no group  $\times$  time interaction effect,

$$H_{01}: (\beta_1 - \beta_4) = (\beta_2 - \beta_5) = (\beta_3 - \beta_6),$$

can be expressed as

$$H_{01}: L\beta = 0,$$

where

$$L = \begin{pmatrix} 1 & -1 & 0 & -1 & 1 & 0 \\ 1 & 0 & -1 & -1 & 0 & 1 \end{pmatrix}.$$

An attractive feature of the general linear model formulation,

$$E(Y_i|X_i) = \mu_i = X_i\beta,$$

is that it can handle settings where the data for some subjects are missing. For example, suppose that the  $i^{th}$  subject belongs to the first group and is missing the response at the third occasion. The appropriate design matrix for that subject is the following  $2 \times 6$  matrix, obtained by removing the last row of the full data design matrix for subjects from the first group:

$$X_i = \begin{pmatrix} 1 & 0 & 0 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 & 0 & 0 \end{pmatrix}.$$

For more general patterns of missingness, the appropriate design matrix for the  $i^{th}$  subject is simply obtained by removing rows of the full data design matrix corresponding to the missing responses. This allows the analysis of response profiles to be based on all available observations of the subjects.

Note that the general linear model for two groups measured at three occasions,

$$E(Y_i|X_i) = \mu_i = X_i\beta,$$

could also have been expressed in terms of the following two design matrices:

$$X_i = \begin{pmatrix} 1 & 0 & 0 & 1 & 0 & 0 \\ 1 & 1 & 0 & 1 & 1 & 0 \\ 1 & 0 & 1 & 1 & 0 & 1 \end{pmatrix},$$

for the first group and

$$X_i = \begin{pmatrix} 1 & 0 & 0 & 0 & 0 & 0 \\ 1 & 1 & 0 & 0 & 0 & 0 \\ 1 & 0 & 1 & 0 & 0 & 0 \end{pmatrix},$$

for the second group. In that case

$$\mu(2) = \begin{pmatrix} \mu_1(2) \\ \mu_2(2) \\ \mu_3(2) \end{pmatrix} = \begin{pmatrix} \beta_1 \\ \beta_1 + \beta_2 \\ \beta_1 + \beta_3 \end{pmatrix};$$

and

$$\mu(1) = \begin{pmatrix} \mu_1(1) \\ \mu_2(1) \\ \mu_3(1) \end{pmatrix} = \begin{pmatrix} \beta_1 + \beta_4 \\ (\beta_1 + \beta_4) + (\beta_2 + \beta_5) \\ (\beta_1 + \beta_4) + (\beta_3 + \beta_6) \end{pmatrix}.$$

The choice of “reference group” (here the second group) is arbitrary, and we have used the convention adopted by many of the procedures in SAS; a more detailed discussion of the “reference group” parameterization is given at the end of this section. With this choice of design matrices for the two groups, the interpretation of the regression coefficients  $\beta$  has changed. Re-expressing hypotheses about the mean response profiles for the two groups in terms of hypotheses about the components of  $\beta$ , we write the hypothesis of no group  $\times$  time interaction as

$$H_{01}: \beta_5 = \beta_6 = 0.$$

Although both alternative parameterizations considered thus far allow for testing of hypotheses about the response profiles, the second parameterization is more convenient since the hypothesis of no group  $\times$  time interaction is represented by the vanishing (or setting to zero) of certain components of  $\beta$ . Also the second parameterization, often called the “reference group” parameterization, is the one that is commonly adopted by many statistical software packages (e.g., PROC MIXED in SAS).

As indicated earlier, when the hypothesis of parallel profiles cannot be rejected, hypotheses concerning the main effects of time and/or group may be of secondary interest, although their relevance depends on the design of the study. Hypotheses concerning the main effects of time and group can similarly be represented by the vanishing (or setting to zero) of certain components of  $\beta$ . For example, with two groups measured at three occasions and assuming parallel profiles ( $\beta_5 = \beta_6 = 0$ ), the hypothesis of no time effect is

$$H_{02}: \beta_2 = \beta_3 = 0;$$

the hypothesis of no group effect is

$$H_{03}: \beta_4 = 0.$$

For the more general case with  $G$  groups measured at  $n$  occasions, the number of constraints under  $H_{02}$  is  $n - 1$  and is the same regardless of the number of groups; the test of  $H_{02}$  has  $n - 1$  degrees of freedom. Similarly the number of constraints under  $H_{03}$  is  $G - 1$  and is the same regardless of the number of occasions; the test of  $H_{03}$  has  $G - 1$  degrees of freedom. Both of these hypotheses can be tested by considering a model with only main effects of group and time. That is, valid tests of the main effects of group and time are obtained from the reduced model that excludes the group  $\times$  time interaction.

Finally, given that the analysis of response profiles can be expressed in terms of the linear regression model

$$E(Y_i|X_i) = \mu_i = X_i\beta,$$

where  $\beta = (\beta_1, \dots, \beta_p)'$  is a  $p \times 1$  vector of regression coefficients (with  $p = G \times n$ ), maximum likelihood estimation of  $\beta$ , and the construction of tests of the group  $\times$  time interaction (and the main effects of time and group), are possible once the covariance of  $Y_i$  has been specified. In the analysis of response profiles, the covariance of  $Y_i$  is usually assumed to be unstructured with no constraints on the  $\frac{n(n+1)}{2}$  covariance parameters other than the requirement that they yield a symmetric matrix that is positive-definite (the condition that the covariance matrix is positive-definite ensures that while the repeated measures can be highly correlated, there must be no redundancy in the sense that one of the repeated measures can be expressed as a linear combination of the others; the condition also ensures that no linear combination of the responses can have a negative variance). Given REML (or ML) estimates of  $\beta$ , and their standard errors (and the estimated covariance of  $\hat{\beta}$ ), tests of the group  $\times$  time interaction (and the main effects of time and group) can be constructed using multivariate Wald tests. Alternatively, likelihood ratio tests can be constructed but require that the model be fit to the data with and without the constraints under the null hypothesis (i.e., fitting the “reduced” and “full” models, respectively). Before illustrating the analysis of response profiles, we present a brief review of the “reference group” parameterization.

# Review: Reference Group Parameterization

Consider a group factor with  $G$  levels. To represent this factor in a linear model, we can define a set of “dummy” or indicator variables:

$$Z_{ig} = \begin{cases} 1 & \text{if } i^{\text{th}} \text{ subject belongs to group } g, \\ 0 & \text{otherwise.} \end{cases}$$

Letting  $X_i = (Z_{i1}, \dots, Z_{iG})$ , the mean response in the  $G$  groups, denoted by  $\mu_i(1), \dots, \mu_i(G)$ , can be expressed in terms of the following linear model:

$$E(Y_i|X_i) = \mu_i = X_i\beta.$$

In this parameterization,

$$\begin{pmatrix} \mu_i(1) \\ \mu_i(2) \\ \vdots \\ \mu_i(G-1) \\ \mu_i(G) \end{pmatrix} = \begin{pmatrix} \beta_1 \\ \beta_2 \\ \vdots \\ \beta_{G-1} \\ \beta_G \end{pmatrix}.$$

If we wish to include an “intercept”, say  $\beta_1$ , by setting the first column of  $X_i$  to 1 (for all  $i = 1, \dots, N$ ), then there is redundancy in  $X_i$  if all  $G$  indicator variables,  $Z_{i1}, \dots, Z_{iG}$ , are also included in the design vector,  $X_i$ . To avoid this over-specification, one of the indicator variables must be excluded from  $X_i$ . Arbitrarily, we can drop  $Z_{iG}$ . Then, with  $X_i = (1, Z_{i1}, \dots, Z_{iG-1})$ , the mean response in the  $G$  groups can be expressed in terms of the following linear model:

$$E(Y_i|X_i) = \mu_i = X_i\beta,$$

where  $\beta = (\beta_1, \dots, \beta G)'$ . In this parameterization,

$$\begin{pmatrix} \mu_i(1) \\ \mu_i(2) \\ \vdots \\ \mu_i(G-1) \\ \mu_i(G) \end{pmatrix} = \begin{pmatrix} \beta_1 + \beta_2 \\ \beta_1 + \beta_3 \\ \vdots \\ \beta_1 + \beta_G \\ \beta_1 \end{pmatrix}.$$

Because the “intercept” term,  $\beta_1$ , is also the mean of group  $G$ , and all of the remaining components of  $\beta$  represent deviations from the mean of group  $G$ , this parameterization is often referred to as the “reference group” parameterization. Here the last level of the group factor (i.e., group  $G$ ) is the reference group, and it is no coincidence that this is the same group whose indicator variable was excluded from  $X_i$ . Other choices of reference group can be obtained by excluding the relevant indicator variable for the group in question from  $X_i$ .

## **5.4 CASE STUDY**

Next we illustrate the main ideas by conducting an analysis of response profiles of the blood lead data of the 100 children from the succimer and placebo groups of the Treatment of Lead-Exposed Children (TLC) Trial.

# Treatment of Lead-Exposed Children Trial

Recall that the TLC trial was a placebo-controlled, randomized trial of an orally administered chelating agent, succimer, in children with confirmed blood lead levels of 20 to 44  $\mu\text{g}/\text{dL}$ . The children in the trial were aged 12 to 33 months and lived in deteriorating inner city housing. The following analysis is based on data on blood lead levels at baseline (or week 0), week 1, week 4, and week 6 during the first treatment period. The mean response profiles for the two groups were displayed in [Figure 5.1](#).

In [Table 5.3](#) the REML estimates of the components of the unstructured covariance matrix are displayed. Note the discernible increase in the variability in blood lead levels from pre- to post-randomization. This increase in variability from baseline is probably due to two factors. First, within each treatment group there may be natural heterogeneity in the individual response trajectories over time. Second, the trial had an inclusion criterion that blood lead levels at baseline were in the range of 20 to 44  $\mu\text{g}/\text{dL}$ ; this may partially account for the smaller variance at baseline.

**Table 5.3** Estimated covariance matrix for the blood lead levels at baseline, week 1, week 4, and week 6 for the children from the TLC trial.

Covariance Matrix			
25.2	19.1	19.7	22.2
19.1	44.3	35.5	29.7
19.7	35.5	47.4	30.6
22.2	29.7	30.6	58.7

In [Table 5.4](#) the results of the analysis of response profiles are presented. The main interest is in the test of the group  $\times$  time interaction. The test of the group  $\times$  time interaction is based on a multivariate Wald test. The test provides a simultaneous test of  $H_0: L\beta = 0$  versus  $H_A: L\beta \neq 0$ , for a suitable choice of  $L$ , and the test statistic can be constructed as

**Table 5.4** Wald tests of fixed effects based on an analysis of response profiles of the blood lead level data at baseline, week 1, week 4, and week 6 for the children from the TLC trial.

Variable	DF	Chi-Squared	P-Value
Group	1	25.43	<0.0001
Week	3	184.48	<0.0001
Group $\times$ Week	3	107.79	<0.0001

$$W^2 = (\hat{L}\hat{\beta})' \{ \hat{L} \hat{\text{Cov}}(\hat{\beta}) \hat{L}' \}^{-1} (\hat{L}\hat{\beta}),$$

and compared to a  $\chi^2$  distribution with degrees of freedom equal to the number of rows of  $L$ . The corresponding likelihood ratio test could be constructed and would require the comparison of the maximized ML log-likelihood for two models, one model that incorporates the constraint that  $L\beta = 0$  (i.e., the model without group  $\times$  time interaction), the other model unconstrained (i.e., the model with group  $\times$  time interaction).

In the TLC trial the question of main scientific interest concerns the comparison of the two treatment groups in terms of their patterns of change from baseline in the mean blood lead levels. This question translates directly into a test of the group  $\times$  time interaction. From [Table 5.4](#) the test of the group  $\times$  time interaction yields a Wald statistic of 107.79 with 3 degrees of freedom (the corresponding likelihood ratio test yields  $G^2 = 74.2$ ). When compared with the reference chi-squared distribution with 3 degrees of freedom, there is strong evidence to reject the null hypothesis and conclude that the patterns of change from baseline are not the same in the two groups. Given the pattern of observed responses (see [Figure 5.1](#)), this result is expected.

The tests of main effects of group and time in [Table 5.4](#) are not meaningful, for two quite different reasons. First, the TLC data are from a randomized trial and the test of the main effect of time is not of subject-matter interest while the test of the main effect of group is subsumed within the test of

group  $\times$  time interaction. Second, in general, the tests of main effects of time and group in [Table 5.4](#) are not meaningful in the presence of a significant group  $\times$  time interaction. This underscores our earlier advice that tests of main effects should only be considered when the assumption of parallel response profiles is tenable. When the profiles are parallel, and there is scientific interest in the main effects of time and/or group, the tests of the main effects of time and group require that the model be re-fit to the data, excluding the group  $\times$  time interaction. The resulting Wald tests for the main effects from this reduced model have the desired interpretations.

So far our analysis of response profiles has provided an omnibus test of the group  $\times$  time interaction. However, unless the test of the group  $\times$  time interaction has only a single degree of freedom, this test does not indicate how the two groups differ. For the latter, we must consider the REML estimates of  $\beta$  and their standard errors presented in [Table 5.5](#); alternative single degree of freedom tests for group  $\times$  time interaction will be discussed in Section 5.5. For ease of interpretation, the baseline (week 0) is chosen as the reference level for time and the placebo group is chosen as the reference level for treatment group. From an examination of the three single-degree-of-freedom contrasts for the group  $\times$  time interaction, the results indicate that children treated with succimer have a discernibly greater decrease in mean blood lead levels from baseline at all occasions when compared to the children treated with placebo. For example, when compared to the placebo group, the succimer group has an additional  $3.152 \mu\text{g}/\text{dL}$  (with SE = 1.257) decrease in mean blood lead levels from baseline to week 6. Of note, there are even larger differences between the two treatment groups earlier in the trial. For example, when compared to the placebo group, the succimer group has an additional  $11.406 \mu\text{g}/\text{dL}$  decrease in mean blood lead levels from baseline to week 1. The apparent rebound in blood lead levels after week 1 in the succimer group is thought to be due to lead that is stored in the bones being mobilized, resulting in a new equilibrium in blood lead levels in the children treated with succimer.

**Table 5.5** Estimated regression coefficients and standard errors based on an analysis of response profiles of the blood lead level data at baseline, week 1, week 4, and week 6 for the children from the TLC trial.

Variable	Group	Week	Estimate	SE	Z
Intercept			26.272	0.710	36.99
Group	S		0.268	1.005	0.27
Week		1	-1.612	0.792	-2.04
Week		4	-2.202	0.815	-2.70
Week		6	-2.626	0.889	-2.96
Group $\times$ Week	S	1	-11.406	1.120	-10.18
Group $\times$ Week	S	4	-8.824	1.153	-7.66
Group $\times$ Week	S	6	-3.152	1.257	-2.51

We remind the reader that our earlier warning about testing main effects in the presence of interactions applies also to the results in [Table 5.5](#). For example, the test of the main effect of group in [Table 5.5](#) ( $Z = 0.27$ ) does not compare the average (over occasions) response in the succimer and placebo groups; instead, it compares the mean response at baseline (here the reference level for time) in the two treatment groups. The lack of equivalence between the tests for the main effect of group in [Tables 5.4](#) and [5.5](#) is a direct consequence of the reference group parameterization adopted here (and commonly used by many statistical software packages; e.g., PROC MIXED in SAS).

Finally, because the TLC data are from a randomized trial, the mean response at baseline is independent of treatment assignment (as was confirmed by the non-significant test of the main effect of group in [Table 5.5](#)). Because of the random assignment to treatment groups, this suggests that the analysis of response profiles could be simplified by fitting a model that omits the main effect of group, thereby forcing the two groups to have the same mean response at baseline. In Sections 5.6 and 5.7 we consider the merits of such an adjustment and compare and contrast alternative strategies for handling the baseline response in different settings. In these sections we highlight how the

analysis of response profiles is a flexible method that can easily be adapted to account for the design of the study and to address questions that are scientifically relevant to any particular study.

## 5.5 ONE-DEGREE-OF-FREEDOM TESTS FOR GROUP BY TIME INTERACTION

As we saw in the previous section, the test for group  $\times$  time interaction is quite general. It posits no specific pattern for the difference in the response profiles between groups. This lack of specificity becomes a problem in studies with a large number of occasions of measurement because the general test for group  $\times$  time interaction, with  $(G - 1) \times (n - 1)$  degrees of freedom, becomes less sensitive to an interaction with a specific pattern as  $n$  increases. Even with as few as three or four occasions of measurement, the general test for interaction will not be as sensitive to specific departures from parallelism as the more focused tests we discuss in this section.

In the typical randomized trial of interventions, subjects are randomized to the intervention groups at baseline and the investigator seeks to determine whether the pattern of response after randomization differs between groups. Randomization implies that the mean at baseline is independent of treatment group; that is, by design, the groups have the same mean response at baseline. In that setting, analysts frequently specify a single contrast believed to best represent the direction in which the pattern of response will differ most markedly. For example, if we assume the first parameterization described in Section 5.3 with two groups and wish to test for equality of the difference between the average response at occasions 2 through  $n$  and the baseline value in the two groups, we can choose the contrast

$$L = (-L_1, L_1),$$

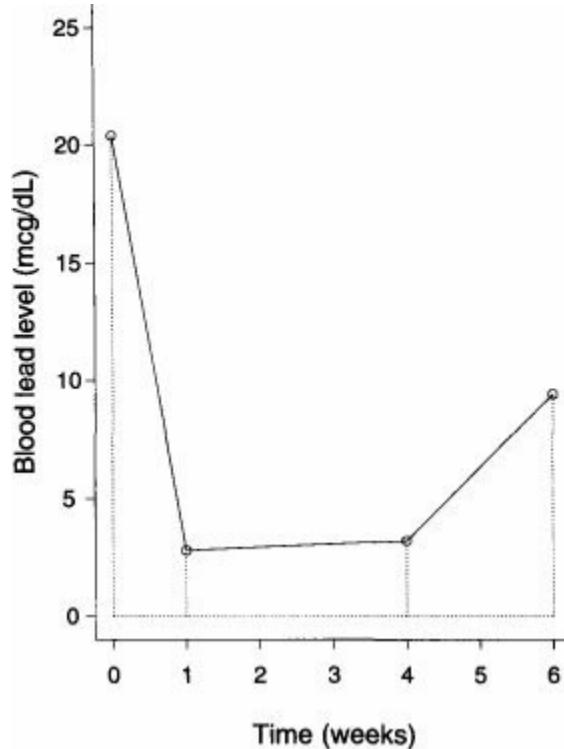
where

$$L_1 = \left( -1, \frac{1}{n-1}, \frac{1}{n-1}, \dots, \frac{1}{n-1} \right).$$

Here  $L_1$  computes the mean response from occasions 2 through  $n$  and subtracts the mean response at baseline for a single group. The latter can be thought of as the average change over the interval for a single group. Thus  $L$  is a group contrast of this average change in the two groups.

A variant of this approach, known as *area under the curve minus baseline*, or sometimes simply AUC, corresponds to a calculation of the area under the trapezoidal curve created by connecting the responses plotted at the respective time points and subtracting  $y_1 x(t_n - t_1)$ , the area of the rectangle of height  $y_1$  and width  $t_n - t_1$ . For illustrative purposes, the AUC of the profile of blood lead levels for a single subject in the TLC trial is shown in [Figure 5.3](#). For this participant, as for most participants in the TLC trial, the AUC is negative because the responses after intervention begins are smaller than the baseline value. The AUC (minus baseline) can be constructed by subtracting the baseline mean,  $\mu_1$ , from each of the means,  $\mu_1$  through  $\mu_n$ , and calculating the area under the trapezoid constructed by connecting these differences. To test for the equality of the AUC in two groups, we would employ the contrast

[\*\*Fig. 5.3\*\*](#) Area under the curve, calculated using the trapezoidal rule, for the profile of blood lead levels for a single subject in the TLC trial.



$$L = (-L_2, L_2),$$

where

$$L_2 = \frac{1}{2} \times (t_1 + t_2 - 2t_n, t_3 - t_1, \dots, t_{j+1} - t_{j-1}, \dots, t_n - t_{n-1})$$

and  $\frac{1}{2} \times (t_{j+1} - t_{j-1})$  is the value of the contrast vector for time points other than 1 (baseline) or  $n$  (the last occasion). These contrast weights are not intuitively obvious, but can be derived from the formula for the area of a trapezoid. Although the curve presented in [Figure 5.3](#) suggests that  $L$  is applied to the individual observations, we must emphasize that the contrast weights are applied to the estimated means, not the individual observations. As a result the AUC can be estimated in settings where some subjects have missing response data.

A third popular method for constructing a single-degree-of-freedom test corresponds to a test of the hypothesis that the trend over time is the same in the several treatment groups. Because this method is a special case of growth curve analysis, to be discussed in depth in the next chapter, and because the expected pattern of response to chelation therapy in the TLC trial would not predict a linear trend in blood lead levels during the treatment period in the group receiving succimer therapy, we defer a discussion of this approach to Chapter 6.

# Application to the Treatment of Lead-Exposed Children Trial

Since the TLC trial measured blood lead levels at four time points during the first treatment period, the vector representing the contrast based on the mean response at times 2 through  $n$  minus baseline is given by

$$L = (-L_1, L_1) = \left(1, -\frac{1}{3}, -\frac{1}{3}, -\frac{1}{3}, -1, \frac{1}{3}, \frac{1}{3}, \frac{1}{3}\right).$$

For the data displayed in [Figure 5.3](#), the change in mean response relative to baseline is  $-15.27$ . Because blood lead levels declined for this subject, as for most participants in the TLC trial, the mean response relative to baseline is negative. From the descriptive statistics in [Table 5.6](#) we can easily determine that the average value of the mean response minus baseline is  $-9.90$  in the succimer group and  $-2.17$  in the placebo group. Thus, if we assume the first parameterization described in Section 5.3, then  $L\hat{\beta} = 7.73$  and the value of the Wald test statistic is  $Z = 8.21$  (or  $W^2 = 67.4$ , with one degree of freedom), indicating a highly significant difference in the response pattern between treatment groups.

**Table 5.6** Mean blood lead levels (and standard deviation) at baseline, week 1, week 4, and week 6 for the children from the TLC trial.

Group	Baseline	Week 1	Week 4	Week 6
Succimer	26.5	13.5	15.5	20.8
	(5.0)	(7.7)	(7.8)	(9.2)
Placebo	26.3	24.7	24.1	23.6
	(5.0)	(5.5)	(5.8)	(5.6)

Similarly, because the time points in the TLC trial were 0, 1, 4, and 6 weeks, the contrast for comparing the AUC (minus baseline) in the two treatment groups is given by

$$L = (-L_2, L_2) = (5.5, -2, -2.5, -1, -5.5, 2, 2.5, 1).$$

The area under the curve for the single subject shown in [Figure 5.3](#) is  $-89.2$ . From [Table 5.6](#), the estimated mean AUC is  $-59.20$  in the succimer group and  $-11.40$  in the placebo group. Thus, if we assume the first parameterization described in Section 5.3, then  $L\hat{\beta} = 47.8$ , yielding a Wald statistic of  $Z = 8.97$  (or  $W^2 = 80.5$ , with one degree of freedom), again highly statistically significant. Thus both methods of analysis provide a clear signal that the response profile differs in the two treatment groups.

Because the TLC trial data provide unequivocal evidence of an effect of succimer on blood lead level, the added sensitivity to treatment effects achieved by the greater specificity of a one-degree-of-freedom test is not important in this application. In many applications, however, the one-degree-of-freedom test will be statistically significant when the overall test for group  $\times$  time interaction is not. For valid application of conventional significance levels, the form of the contrast must be specified prior to data analysis. Otherwise, one would be at risk of seeking the best contrast and testing its significance as if it had been chosen in advance. To guard against this criticism, the protocols for randomized trials usually specify the form of the contrast. This requirement highlights a hazard of one-degree-of-freedom tests. The added sensitivity comes at the price of reduced generality. If the difference between treatment groups takes a form quite different from the pattern anticipated by the contrast, one can fail to obtain a statistically significant result for a one-degree-of-freedom test even when the overall test for group  $\times$  time interaction is statistically significant. Thus one-degree-of-freedom tests should be employed only when there is sufficient prior information to specify the contrast with confidence.

## 5.6 ADJUSTMENT FOR BASELINE RESPONSE

When the data are complete, each of the one-degree-of-freedom tests described in the previous section can be constructed by calculating a univariate summary statistic for each study participant and performing a test for equality of means of these summary statistics in the  $G$  groups. With complete data, group comparisons of these summary statistics are equivalent to applying the corresponding contrast weights to the mean responses. This is because the difference in the means is the mean of the differences when each subject is measured at every occasion. Moreover, for each of the two tests described in detail, mean response minus baseline and AUC minus baseline, the summary statistic corresponds to subtracting the baseline value from a summary of the responses on occasions 2 through  $n$ . For example, for the test for equality of mean response minus baseline, the summary statistic for the  $i^{th}$  participant is given by

$$(5.1) \quad \frac{(Y_{i2} + Y_{i3} + \dots + Y_{in})}{n - 1} - Y_{i1}.$$

With this representation in mind, some analysts have suggested an alternative approach analogous to analysis of covariance (ANCOVA), in which a summary of the response at times 2 through  $n$  becomes the dependent variable and the baseline value enters the analysis as a covariate. When the response variable is the mean at occasions 2 through  $n$  and we wish to test for the equality of the mean in two treatment groups, we can write the corresponding univariate model as

$$(5.2) \quad Y_i^* = \beta_1 + \beta_2 Y_{i1} + \beta_3 \text{trt}_i + e_i^*,$$

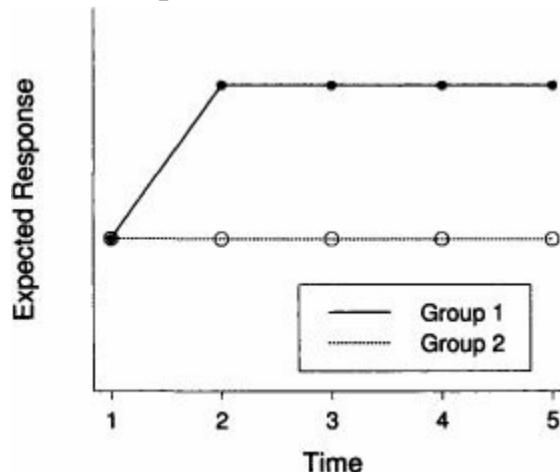
where

$$Y_i^* = \frac{(Y_{i2} + Y_{i3} + \dots + Y_{in})}{n - 1}$$

is the mean response at occasions 2 through  $n$  for the  $i^{th}$  subject,  $\text{trt}_i$  is an indicator variable distinguishing the two treatment groups, and  $e_i^*$  is the error term in the univariate model. This model assumes that the data are complete, and it cannot be fit with missing data; we defer a discussion of more general approaches for handling baseline response to Section 5.7.

An analysis based on either (5.1) or (5.2) will be especially appealing in settings where initial changes from baseline are expected to persist throughout the duration of follow-up. For example, in a trial where the impact of the intervention on changes in the mean response at the start of follow-up is expected to be similar to that toward the end of follow-up; this pattern for the mean response profiles is illustrated in [Figure 5.4](#). Tests based on (5.1) or (5.2) correspond to a comparison between groups of the mean responses on occasions 2 through  $n$ , with adjustment for baseline, and have  $G - 1$  degrees of freedom irrespective of the number of occasions of measurement.

**Fig. 5.4** Graphical representation of changes in the mean response from baseline (in Group 1) that persist throughout the duration of follow-up.



This raises a question about whether one should incorporate the baseline value through the contrast given by (5.1) or through the analysis of covariance model given by (5.2) in a specific application. The answer depends critically on whether the data arose from an observational study or a

randomized trial. If the study is an observational one, for example, a longitudinal study of the determinants of rate of decline of pulmonary function in adults, it is usually not advisable to employ the analysis of covariance approach because the baseline value may be associated with other variables whose effects are to be studied, raising problems of confounding in an analysis intended to describe how the pattern of response over time is influenced by the characteristics of study participants. For example, individuals who are smokers as adults might have smoked during adolescence. If smoking affected the attained pulmonary function level for young adults, then smoking will likely be associated with pulmonary function level later in adult life, even if cigarette smoking does not influence the rate of decline of pulmonary function with age. Thus adjustment for baseline pulmonary function level using (5.2) could introduce an association between smoking status and rate of decline of pulmonary function, even if the unadjusted rates of decline are nearly equal in the various smoking groups.

When participants have been randomized to the several treatment groups and the baseline value has been obtained before any study interventions, adjustment for baseline through analysis of covariance is of interest. In that setting the mean response at time 1 is independent of treatment assignment. One can then show that the one-degree-of-freedom test for equality of response profiles based on a contrast and the corresponding test based on analysis of covariance represent alternative tests of the same null hypothesis and that the test based on the analysis of covariance approach will always be more efficient. That is, the analysis of covariance approach yields estimates of treatment effects with smaller standard errors than those obtained by calculating contrasts.

For example, the greater efficiency of the analysis of covariance can be highlighted by examining the relative efficiency of (5.1) to (5.2) in simple settings. The relative efficiency is defined as the ratio of the variance of the estimator based on (5.2) to the variance of the estimator based on (5.1); a relative efficiency less than one implies that the estimator based on (5.2) has smaller variance. When the covariance among the repeated measures is assumed to have a compound symmetry pattern, with common variance  $\sigma^2$  and common correlation  $\rho$ , the relative efficiency is given by

$$(5.3) \quad \frac{1}{n} \{1 + (n - 1)\rho\}.$$

The derivation of (5.3) is not important. What this simple expression indicates is that the two methods of adjustment for baseline response are equally efficient only when  $\rho = 1$ . When  $\rho = 0$ , the analysis based on (5.1) is only  $1/n$  times as efficient as the analysis of covariance. The greater efficiency of the analysis of covariance depends on both the number of repeated measures and the strength of the correlation among them. For example, when  $n = 5$  and  $\rho = 0.4$ , the analysis of covariance is approximately twice as efficient as subtracting the baseline response.

In general, the analysis of longitudinal data from a randomized trial is the only setting where we recommend adjustment for baseline through analysis of covariance. In that setting, in contrast to observational studies, adjustment leads to meaningful tests of hypothesis of scientific interest. Moreover the tests based on the analysis of covariance approach will be more powerful. The notion of adjustment for baseline can also be applied more generally in the analysis of response profiles; in Section 5.7 we compare and contrast a number of alternative strategies for handling the baseline response in more general settings and make recommendations about the preferred strategies in different situations.

We conclude this section by noting that adjustment for baseline in the analysis of longitudinal change is a topic that has generated heated debate among analysts. When longitudinal data arise from an observational study, the two methods of adjusting for baseline described in this section can yield discernibly different and, apparently conflicting, results. This conundrum is also known as *Lord's paradox* (named after Frederic Lord, who eloquently brought the issue to light) and has led many researchers astray over the years. The paradox lies in the interpretation of the two types of analyses and is resolved by noting that these two alternative methods of adjusting for baseline answer qualitatively different scientific questions when the data arise from an observational study. This can be illustrated in the simplest setting where there are two groups or sub-populations (e.g., males and females) measured at two occasions. The overall goal of such a study is to compare the changes in

response for the two groups. The analysis that subtracts baseline response, thereby creating a simple change score, addresses the question of whether the two groups differ in terms of their mean change over time. In contrast, adjustment for baseline using analysis of covariance addresses the question of whether an individual belonging to one group is expected to change more (or less) than an individual belonging to the other group, *given that they have the same baseline response*. The latter question is a conditional one and, depending on the study design, may address a different scientific question than the former.

For example, in an observational study examining gender difference in weight gain of infants between the ages of 12 and 24 months, a measure of body weight might be obtained at 12 months (baseline) and at 24 months. The analysis of the simple change score addresses the question of whether boys and girls differ in terms of their changes in mean body weight over the 12 months of follow-up. At baseline, boys are on average  $1\frac{1}{2}$  pounds heavier than girls, but there is no evidence of a gender effect on the 12 month changes in body weight, with boys and girls both gaining approximately  $5\frac{1}{4}$  pounds. In contrast, the analysis of covariance of the same data reveals a discernible gender effect, with boys showing more weight gain than girls. Thus, even though the unadjusted (or *unconditional*) increases in body weight are approximately the same for this age cohort of boys and girls, the analysis of covariance is directed at the *conditional* question of whether boys are expected to gain more weight than girls, *given that they have the same initial weight at 12 months*. That is, if we compare boys and girls within sub-populations with the same initial weight at 12 months, are their average weights at 24 months the same? When the conditional question is posed this way, we would expect boys to gain more weight than girls. The reasoning is that if a boy and girl have the same initial weight at 12 months, then there are two possibilities: (1) the girl is initially overweight and is expected to gain less weight over the 12 months, or (2) the boy is initially underweight and is expected to gain more.

A more thorough discussion of this issue is beyond the scope of this book, but we advise readers to employ the analysis of covariance approach in longitudinal settings only if the approach and its implications are fully understood.

In summary, the choice between the two methods of adjusting for baseline discussed in this section should be made on substantive grounds. That is, the design of the longitudinal study and the research question of interest should guide the choice of analytic method. The analysis that subtracts baseline response is appropriate when the primary goal of the study is to compare distinct populations in terms of their average change over time. The analysis addresses the question: Do the populations differ in terms of their average change? and this is appropriate when the data have arisen from either an observational study or a randomized trial. On the other hand, analysis of covariance will, in general, be appropriate in cases where individuals have been assigned to groups at random (e.g., a randomized trial) or where the population distributions of the baseline responses can reasonably be assumed to be equal (even though the sample means of the baseline responses may differ across groups). In cases where the population distributions of the baseline responses are equal, it is then meaningful to ask the question: Is the expected change the same in all groups, when we compare individuals having the same baseline response? Furthermore the analysis of covariance will provide a more powerful test of group differences. The latter has often been touted as the main reason that analysis of covariance should be the preferred method of adjusting for baseline. This faulty rationale, however, has blinded many researchers to the potential difficulties in interpreting the results of analysis of covariance when the assumption of equal population distributions of baseline response is not tenable. In conclusion, it is the study design and the scientific question of interest, and not issues of statistical precision and power, that should primarily determine the choice of analytic methods for adjusting for baseline response.

## 5.7 ALTERNATIVE METHODS OF ADJUSTING FOR BASELINE RESPONSE\*

One feature of longitudinal studies that sets them apart from repeated measures and related designs is the presence of a baseline measurement. In a randomized longitudinal trial comparing treatments, the measurement at the first occasion is usually a baseline response obtained prior to any study interventions. For example, in the TLC trial, the blood lead levels at baseline were obtained prior to receiving placebo or succimer. In that case, due to randomization, we can assume that the treatment group means are equal at baseline. Thus the question naturally arises as to how to handle the baseline measurement in the assessment of whether patterns of change in the mean response over time are the same in the groups.

In the previous section we considered two methods of adjustment for baseline in a relatively simple setting. The notion of adjustment for baseline can also be applied more generally in the analysis of response profiles. In this section<sup>†</sup> we compare and contrast a number of alternative strategies for handling the baseline response and make recommendations about the preferred strategies in different situations.

We consider four ways of handling the baseline value:

1. We can retain it as part of the outcome vector and make no assumptions about group differences in the mean response at baseline.
2. We can retain it as part of the outcome vector and assume the group means are equal at baseline, as might be appropriate in a randomized trial.
3. We can subtract the baseline response from all of the remaining post-baseline responses, and analyze the differences from baseline.
4. We can use the baseline value as a covariate in the analysis of the post-baseline responses.

We now consider the appropriateness and merits of each of these four strategies and illustrate their application to the blood lead level data from the TLC trial.

The first two strategies retain the baseline measurement as part of the outcome vector, but differ in terms of assumptions about the mean response at baseline. The first strategy corresponds to a standard analysis of response profiles without incorporating any constraints on the group means at baseline. This was the method of analysis highlighted in Sections 5.2 through 5.4. The second strategy corresponds to an analysis of response profiles where the group means at baseline are constrained to be equal. In a randomized trial, where treatment assignment is random, both strategies yield valid estimates of treatment group comparisons, but the second strategy is, in general, more powerful. Thus, in randomized trials, or in observational studies where there is good reason to assume that the groups have the same mean response at baseline (e.g., due to matching on baseline response), the second strategy for handling baseline is preferred and should be routinely used. In contrast, in observational studies where there is no a priori reason to assume the groups have the same mean response at baseline, the second strategy is not appropriate and only the first strategy should be used.

In the analyses of the blood lead level data from the TLC trial presented in Section 5.4, the first strategy was employed in the results presented in [Tables 5.4](#) and [5.5](#). The second strategy can be implemented by excluding the treatment group main effect from the model for the response profiles. This model is unusual in that it contains an interaction between group and time but no main effect of group. This model appears to contradict the conventional wisdom that interactions should not be included in a regression model without their main effects. However, this is an important exception to the rule. Because baseline (week 0) was chosen as the reference level for time, the exclusion of the group main effect forces the two groups to have the same mean response at baseline. The results of such an analysis are presented in [Table 5.7](#) and are qualitatively similar to those in [Table 5.5](#). Note the absence of the main effect of group, which has been set to zero. Also the omnibus test of the group  $\times$  time interaction from this model yields a Wald statistic of 111.96 with 3 degrees of freedom. In

contrast, the analysis of response profiles without any adjustment for baseline (strategy 1) yielded a Wald statistic of 107.79, with 3 degrees of freedom (see [Table 5.4](#)). The difference between these two statistics reflects the increased power of the second method for handling baseline response.

**Table 5.7** Estimated regression coefficients and standard errors based on an analysis of response profiles of the blood lead level data assuming equal mean blood lead levels at baseline in the succimer and placebo groups.

Variable	Group	Week	Estimate	SE	Z
Intercept			26.406	0.500	52.83
Week		1	-1.645	0.782	-2.10
Week		4	-2.231	0.807	-2.76
Week		6	-2.642	0.887	-2.98
Group × Week	S	1	-11.341	1.093	-10.38
Group × Week	S	4	-8.765	1.131	-7.75
Group × Week	S	6	-3.120	1.251	-2.49

The third and fourth strategies do not retain the baseline response as part of the outcome vector. Instead, they focus on raw and adjusted changes from baseline and restrict the outcome vector to measurements obtained post-baseline. The third strategy is to subtract the baseline response from the remaining post-baseline responses, and analyze the differences from baseline. We refer to these differences from baseline as “raw change scores.” With responses at  $n$  occasions, we can define the  $(n - 1) \times 1$  vector of raw change scores,

$$D_i = (Y_{i2} - Y_{i1}, Y_{i3} - Y_{i1}, \dots, Y_{in} - Y_{i1})'$$

and conduct an analysis of response profiles with  $D_i$  as the outcome vector. Because the outcome is a change score (the change from baseline), this approach alters the interpretation of the tests for all three effects in the analysis of response profiles. The test for group  $\times$  time interaction becomes a test for parallel profiles for the changes from baseline in the mean response on occasions 2 through  $n$ ; the test for group effect becomes a test that the changes from baseline at occasion 2 are the same across groups (assuming that occasion 2 is chosen as the reference level for time). Thus, to address the question of whether patterns of change over time are the same in all groups, the test of interest under the third strategy must be modified. It is now a joint test that combines the main effect of group and the group  $\times$  time interaction. The test has the same  $(G - 1) \times (n - 1)$  degrees of freedom because it combines a  $(G - 1)$  degrees of freedom test for group with a  $(G - 1) \times (n - 2)$  degrees of freedom test for group  $\times$  time interaction. This is in contrast to the conventional analysis of response profiles (with baseline response included as part of the response vector) where only the group  $\times$  time interaction addresses important questions concerning group comparisons of the patterns of changes in the mean response and the group main effect is not of scientific interest.

Of note, this joint test of the main effect of group and the group  $\times$  time interaction is formally equivalent to the test of the group  $\times$  time interaction (with  $(G - 1) \times (n - 1)$  degrees of freedom) under the first strategy. Moreover, for the purposes of analyzing changes in the mean response and how these changes differ among groups, the first and third strategies are completely equivalent; that is, the first and third strategies for handling baseline produce identical tests and estimates of effects. Thus it is clear that the third strategy offers no efficiency gain.

Using the blood lead level data from the TLC trial, the results of the analysis of change scores are presented in [Table 5.8](#). At first glance the regression parameter estimates in [Table 5.8](#) do not appear to agree with those in [Table 5.5](#). However, it can easily be shown that the six parameter estimates in [Table 5.8](#) are simple linear combinations of the estimates for the time and group  $\times$  time interaction effects reported in [Table 5.5](#). For example, the group effect, -11.406, in [Table 5.8](#) is identical to the estimate of the group  $\times$  time interaction in [Table 5.5](#) that represents the group contrast of the changes from baseline (week 0) to week 1. Similarly the estimated group contrast of the changes from baseline to week 4, -8.824, from [Table 5.5](#) can also be obtained from the estimates in [Table 5.8](#) ( $-11.406 + 2.582 = -8.824$ ). It can be shown that all estimates of change from baseline, and the

group comparisons of these changes, reported in [Tables 5.5](#) and [5.8](#) are identical.

**Table 5.8** Estimated regression coefficients and standard errors based on an analysis of response profiles of the changes from baseline in blood lead levels at week 1, week 4, and week 6 for the children from the TLC trial.

Variable	Group	Week	Estimate	SE	Z
Intercept			-1.612	0.792	-2.04
Group	S		-11.406	1.120	-10.18
Week		4	-0.590	0.643	-0.92
Week		6	-1.014	0.934	-1.09
Group × Week	S	4	2.582	0.909	2.84
Group × Week	S	6	8.254	1.321	6.25

The analysis of change scores reported in [Table 5.8](#) produced a Wald statistic of 107.79, with 3 degrees of freedom, for jointly testing the effect of group and the group  $\times$  time interaction. As expected, this agrees with the omnibus test of the group  $\times$  time interaction reported in [Table 5.4](#) from the analysis of response profiles without any adjustment for baseline (strategy 1).

Because the first and third strategies yield identical analyses of changes in the mean response, and how these changes differ among groups, in principle either method can be used. However, from a practical standpoint, we recommend the first strategy over the third for two main reasons. First, the analysis of change scores has implications for the interpretation of the hypothesis tests that are more consequential. For example, while the test of primary interest in the conventional analysis of response profiles is usually the test for group  $\times$  time interaction (with  $(G - 1) \times (n - 1)$  degrees of freedom), the test of interest in the analysis of change scores is a  $(G - 1) \times (n - 1)$  degrees of freedom test that incorporates the main effect of group (with  $(G - 1)$  degrees of freedom) and the  $(G - 1) \times (n - 2)$  degrees of freedom group  $\times$  time interaction in this model. The analysis of change scores requires the construction of joint tests of main effects and interactions; these tests are not routinely produced as standard output from statistical software for analyzing response profiles. Second, when there are subjects with missing baseline response, all their data are excluded from the analysis of change scores; in contrast, the first strategy incorporates all available data in the analysis.

The fourth strategy is to analyze the post-baseline responses and make an adjustment for the baseline response by including it as a covariate. If we have responses at  $n$  occasions, the analysis of response profiles is based on the  $(n - 1) \times 1$  vector,

$$Y_i = (Y_{i2}, Y_{i3}, \dots, Y_{in})'$$

and the baseline response,  $Y_{i1}$ , is regarded as a covariate. This type of analysis corresponds to an analysis of covariance (ANCOVA), albeit one where the outcome is a  $(n - 1) \times 1$  vector of responses. This strategy for handling baseline is appropriate when analyzing data from randomized trials. Note that because of randomization, hypotheses of equality of the *conditional* means of the response at occasions 2 through  $n$ , given the baseline response, imply hypotheses of equality of the *unconditional* means of the response at occasions 2 through  $n$ . This strategy for handling baseline is appropriate also for observational studies where there is good reason to assume the groups have the same mean response at baseline. It should not be used, however, in observational studies where there is no a priori reason to assume that the groups have the same mean response at baseline.

Interestingly, the fourth strategy for handling baseline can be implemented by conducting the analysis of response profiles (with  $Y_{i1}$  included as a covariate) on either the post-baseline responses or the post-baseline change scores. That is, the estimates of all effects of interest are identical whether the analysis is based on the  $(n - 1) \times 1$  vector of post-baseline response,  $Y_i = (Y_{i2}, Y_{i3}, \dots, Y_{in})'$ , or the  $(n - 1) \times 1$  vector of post-baseline differences,  $D_i = (Y_{i2} - Y_{i1}, Y_{i3} - Y_{i1}, \dots, Y_{in} - Y_{i1})'$ . The intuition for why these two analyses are identical is as follows: Because the two outcomes differ by  $Y_{i1}$ , and both analyses estimate effects that are adjusted for  $Y_{i1}$  by holding the baseline value fixed, they produce the same regression coefficients for all effects of interest. The two analyses differ only

in terms of the estimated slope for  $Y_{i1}$ . (The estimated slope from the analysis based on  $Y_i$  is simply one unit larger than the estimated slope from the analysis based on  $D_i$ .) Because of this equivalence, we can regard the fourth strategy as an analysis of the “adjusted change scores” (i.e.,  $D_i$  adjusted for  $Y_{i1}$ ) in contrast to an analysis of the raw or unadjusted change scores (strategy 3).

Because the outcome is an adjusted change score, the fourth strategy for handling baseline also alters the interpretation of the tests for all three effects in the analysis of response profiles. The test for group  $\times$  time interaction becomes a test for parallel profiles for the adjusted changes from baseline in the mean response on occasions 2 through  $n$ ; the test for group effect becomes a test that the adjusted changes from baseline in the mean response at occasion 2 are the same across groups (assuming that occasion 2 is chosen as the reference level for time). Thus, similar to the third strategy, the test of interest is a joint test of the main effect of group and the group  $\times$  time interaction.

Using the blood lead level data from the TLC trial, the results of the analysis of adjusted change scores are presented in [Table 5.9](#). Although the results of this analysis are qualitatively similar to the results from the analysis of raw change scores (strategy 3), note that the estimated main effect of group (and the intercept) in [Table 5.9](#) is slightly different from the corresponding estimate in [Table 5.8](#).

**Table 5.9** Estimated regression coefficients and standard errors based on an analysis of response profiles of the adjusted changes from baseline in blood lead levels at week 1, week 4, and week 6 for the children from the TLC trial.

Variable	Group	Week	Estimate	SE	Z
Intercept			-1.638	0.777	-2.11
Baseline <sup>a</sup> ( $Y_{i1} - 26.406$ )			-0.196	0.094	-2.08
Group	S		-11.354	1.099	-10.34
Week		4	-0.590	0.643	-0.92
Week		6	-1.014	0.934	-1.09
Group $\times$ Week	S	4	2.582	0.909	2.84
Group $\times$ Week	S	6	8.254	1.321	6.25

<sup>a</sup>Centering baseline response on its overall mean (26.406) gives the intercept a meaningful interpretation.

The analysis of adjusted change scores reported in [Table 5.9](#) produced a Wald statistic of 111.13, with 3 degrees of freedom, for jointly testing the effect of group and the group  $\times$  time interaction. This is larger than the corresponding statistic under the third strategy, reflecting the greater efficiency of the analysis of adjusted change scores. The greater efficiency of analysis of covariance (adjusted change score analysis) over simple contrasts (raw change score analysis) was highlighted in Section 5.6.

Observe that the Wald statistic of 111.13 produced by the analysis of adjusted change scores is quite similar to that obtained under the second strategy for adjusting for baseline. Thus the question naturally arises as to which approach is preferred: strategy 2 or strategy 4. One could argue that strategy 2 is preferred over strategy 4 for exactly the same reasons given for preferring strategy 1 over strategy 3. That is, the analysis of adjusted change scores requires the construction of joint tests of main effects and interactions, and these tests are not routinely produced as standard output from statistical software for analyzing response profiles. Second, when there are subjects with missing baseline response, all their data are excluded from the analysis of adjusted change scores; in contrast, the second strategy incorporates all available data in the analysis. Finally, there is a third reason why strategy 2 might be preferred over strategy 4. There is an implicit assumption in the adjusted change score analysis that the regression slope relating  $Y_{ij}$  to  $Y_{i1}$  (for  $j = 2, \dots, n$ ) is the same at all  $n - 1$  post-baseline occasions. This implies a strong assumption about the covariances among  $Y_{i1}, \dots, Y_{in}$ . Specifically, it constrains

$$\text{Cov}(Y_{i1}, Y_{i2}) = \text{Cov}(Y_{i1}, Y_{i3}) = \dots = \text{Cov}(Y_{i1}, Y_{in}).$$

As a result of these constraints, there is the potential for misspecification of the model for the

covariance and misleading inferences about change over time. In contrast, the second strategy imposes no such structure on the covariance matrix. Finally, we note that if the assumption of homogeneous regression slopes (relating  $Y_{ij}$  to  $Y_{i1}$ , for  $j = 2, \dots, n$ ) is relaxed and  $n - 1$  separate regression slopes are estimated, then strategies 2 and 4 yield identical estimates of effects, but strategy 2 yields slightly more powerful tests. Thus the second strategy for handling baseline can be seen to enjoy all of the efficiency gains that have been highlighted in Section 5.6 for ANCOVA (adjusted change score analysis). Therefore, on practical grounds, for the reasons outlined above, it can be argued that strategy 2 is preferred over strategy 4.

To summarize, there are many ways to handle baseline response in the analysis of longitudinal data. In this section we have reviewed four strategies. We have seen that the two methods that retain the baseline value as part of the outcome are completely equivalent to corresponding strategies that restrict the outcome vector to measurements obtained post-baseline. However, on practical grounds, it can be argued that the first and second strategies are preferable. The first and second strategies differ in terms of efficiency. The second strategy enjoys all the efficiency gains that have been highlighted in Section 5.6 for ANCOVA. But the choice between the first and second strategies should be guided by the study design. In randomized trials, or in observational studies where there is good reason to assume the groups have the same mean response at baseline, the second strategy is, in general, more powerful and should be routinely used. In contrast, in observational studies where there is no a priori reason to assume that the groups have the same mean response at baseline, the second strategy is not appropriate and the first should be used.

## 5.8 STRENGTHS AND WEAKNESSES OF ANALYZING RESPONSE PROFILES

The analysis of response profiles is a conceptually straightforward way to analyze data from a longitudinal study when the design is balanced, with the timing of the repeated measures common to all individuals in the study, and when all the covariates are discrete (e.g., representing different treatments, interventions, or characteristics of the study subjects). The main feature of the analysis of response profiles is that it allows arbitrary patterns in the mean response over time and arbitrary patterns in the covariance of the responses. As a result this method for longitudinal analysis has a certain robustness, since the potential risks of bias due to misspecification of the models for the mean and covariance are minimal. Although the analysis of response profiles requires that the data arise from a balanced design, it can be applied when the data are incomplete due to missing response data.

The method for analyzing response profiles described in this chapter is related to a more traditional approach known in the statistical literature as “profile analysis”. However, we make a distinction between the method presented in this chapter and a traditional profile analysis. In a traditional profile analysis the three hypotheses concerning response profiles described at the beginning of Section 5.2 are placed on an equal footing, and there is an overwhelming emphasis on hypothesis testing rather than estimation of effects. As we have seen, however, tests of hypotheses concerning main effects of time and/or group often have no direct bearing on questions of scientific interest in a longitudinal study. This is especially the case for longitudinal data arising from randomized trials. Consequently the routine use of traditional profile analysis for longitudinal data coerces the analyst to test certain hypotheses that do not necessarily translate into meaningful scientific questions about longitudinal change in the response. In addition, because profile analysis is often implemented within a multivariate analysis of variance (MANOVA) that requires a complete response vector on each individual, it does not permit subjects with missing responses. The resulting analysis is very inefficient because it is based only on data from the so-called complete cases; it can also produce biased estimates of change in the mean response over time when such “completers” are not a random sample from the target population. Finally, traditional profile analysis lacks flexibility in handling the baseline response and requires that it be part of the response vector. In contrast, the method for analyzing response profiles presented in this chapter can be readily adapted to address specific questions that are well grounded in the science, can be applied when the data are incomplete due to missing response data, and permits alternative approaches for making adjustments for the baseline response.

Although it was not considered here, the analysis of response profiles can be extended in a straightforward way to handle the case where individuals can be grouped according to more than a single factor. For example, if there are two covariates that are discrete (e.g., treatment group and gender), the analysis will include tests of the 3-way and 2-way interactions among these two factors and time (in addition to their main effects). The general linear model can also be used to provide estimated means for summary measure analyses that are based on linear combinations of the mean response vector, for instance, “area under the curve” analysis, when the data are incomplete.

The analysis of response profiles does have a number of potential drawbacks that make it either unappealing or unsuitable for analyzing data from many longitudinal studies. First, the requirement that the longitudinal design be balanced implies that the method cannot be applied when the vectors of repeated measures are obtained at different sequences of time, except by “moving” an observation to the nearest planned measurement time. As a result the method is not well suited to handle mistimed measurements, a common problem in many longitudinal studies. Note, however, that the general method for analyzing response profiles can handle unbalanced patterns of observations due to missing response data. Second, the analysis of response profiles ignores the time ordering of the repeated measures in a longitudinal study. Indeed, the analysis of response profiles could be applied when each individual has a vector of multivariate outcomes that are distinct and non-commensurate (i.e., measures of more than one outcome) rather than repeated measures of a single outcome.

Because the analysis of longitudinal response profiles allows for an arbitrary pattern in the mean responses, and does not impose any time trends, the results of the analysis provide only a very broad or general statement about group differences in patterns of change over time. Ordinarily a significant group  $\times$  time interaction effect that has more than a single degree of freedom will require additional analysis to provide a more informative description of how the groups differ in their patterns of change in the mean response. Third, because the analysis of response profiles produces an overall or omnibus test of effects, it may have low power to detect group differences in specific trends in the mean response over time (e.g., linear trends in the mean response). Single-degree-of-freedom tests of specific time trends are more powerful. Finally, in the analysis of response profiles, the number of estimated parameters ( $G \times n$  mean parameters and  $\frac{n(n+1)}{2}$  covariance parameters) grows rapidly with the number of measurement occasions. For example, with two groups measured at three occasions, the number of parameters is 12. However, with two groups measured at 10 occasions, the number of parameters is 75. Consequently this method is more appealing when the total number of subjects,  $N$ , is relatively large in comparison to the number of measurement occasions,  $n$ .

# 5.9 COMPUTING: ANALYZING RESPONSE PROFILES USING PROC MIXED IN SAS

The MIXED procedure in SAS is a very general and versatile procedure for fitting linear models to longitudinal and clustered data. No attempt is made here to give a comprehensive review of the main features of PROC MIXED. Instead, we present illustrative source code for an analysis of response profiles in general terms and then describe the most salient parts of the command syntax. Many of the later chapters will include a description of additional commands and features of PROC MIXED as they are needed. Although these concluding sections in each chapter will not provide a training manual for the use of PROC MIXED, they should provide a firm basis for understanding the command syntax required for analyzing longitudinal data using PROC MIXED in SAS.

Before discussing the command syntax for PROC MIXED, we note that the procedure requires each repeated measurement in a longitudinal data set to be a separate “record.” For example, in the TLC trial, the data are recorded as follows:

ID	Group	Baseline	Week 1	Week 4	Week 6
001	P	30.8	26.9	25.8	23.8
002	S	26.5	14.8	19.5	21.0
003	S	25.8	23.0	19.1	23.2
004	P	24.7	24.5	22.0	22.5
005	S	20.4	2.8	3.2	9.4
006	S	20.4	5.4	4.5	11.9
:	:	:	:	:	:
100	P	31.1	31.2	29.2	30.1

with a single “record” of the four repeated measurements for each child in the study. When the data set is in this form, it is said to be in a *multivariate* mode or *wide* format. Prior to analysis these data must be converted to a data set with four records for each child, one for each measurement occasion. In the latter form the data set is said to be in a *univariate* mode or *long* format. This can be accomplished using the illustrative SAS commands in [Table 5.10](#), which produce the following data set:

ID	Group	Time	Y
001	P	0	30.8
001	P	1	26.9
001	P	4	25.8
001	P	6	23.8
002	S	0	26.5
002	S	1	14.8
002	S	4	19.5
002	S	6	21.0
:	:	:	:
100	P	0	31.1
100	P	1	31.2
100	P	4	29.2
100	P	6	30.1

[Table 5.10](#) Illustrative commands in SAS for transforming a data set with a single record for each individual to a data set with multiple records corresponding to each measurement occasion.

DATA lead;

```
INFILE 'tlc.dat';
INPUT id group $ y1 y2 y3 y4;
```

```
y=y1; time=0; OUTPUT;  
y=y2; time=1; OUTPUT;  
y=y3; time=4; OUTPUT;  
y=y4; time=6; OUTPUT;  
DROP y1-y4;
```

---

To conduct an analysis of response profiles with data from two or more treatment groups measured repeatedly over time, we can use the illustrative SAS commands given in [Table 5.11](#). This model assumes that the covariance matrix is unstructured. Alternative assumption about the covariance can be considered, and this may be advantageous when the number of measurement occasions is relatively large in comparison to the number of subjects. Choosing a model for the covariance matrix is a topic that will be discussed in Chapter 7. Next we present a brief description of each of the command statements in [Table 5.11](#).

**Table 5.11** Illustrative commands for an analysis of response profiles using PROC MIXED in SAS.

---

```
PROC MIXED;  
  CLASS id group time;  
  MODEL y=group time group*time/S CHISQ;  
  REPEATED time / TYPE=UN SUBJECT=id R RCORR;
```

---

PROC MIXED <options>;

The PROC MIXED statement calls the procedure MIXED in SAS. It can also include an option for the choice of method of estimation. By default, PROC MIXED uses REML estimation; ML estimation can be invoked by including the option METHOD=ML.

CLASS variables;

The CLASS statement is used to define all variables that are to be regarded as categorical or factors. By default, this statement will create indicator variables for each factor using a reference group coding, with the last level (where “last” here refers to the level with the largest alpha-numeric value) regarded as the reference group. Of note, this default indicator variable coding is not the most natural for a categorical variable denoting the occasions of measurement; for the latter, the “first” level of the factor (e.g., the baseline measurement occasions) is usually the natural reference group.

The default coding can be changed with the inclusion of the ORDER= option in the PROC MIXED statement. For example, the ORDER=DATA option forces the levels of all variables included in the CLASS statement to be sorted by their order of appearance in the input data set. Therefore, by previously sorting the data set in descending order of a categorical variable denoting the occasions of measurement, a more natural reference group coding is obtained for that variable (with the lowest, rather than the highest, level of the categorical variable for time used as the reference). However, one unappealing consequence of circumventing the default coding of time in this way is that the estimates of the covariance matrix are printed in reverse (or descending) order of time. To avoid potential confusion when extracting the estimates of the covariance matrix, it is advisable to re-run the analysis without the ORDER=DATA option.

MODEL dependent = <fixed effects> / <options>;

The MODEL statement specifies the response variable and the fixed effects. The fixed effects can include both discrete (defined in the CLASS statement) and quantitative (excluded from the CLASS statement) covariates.

The covariates included in the MODEL statement determine the design matrix  $X_i$ . Of note, by default, PROC MIXED includes a column of 1's in  $X_i$  for the intercept. The option NOINT requests that no intercept be included in the model.

Various options that can be included on the MODEL statement modify how test statistics are computed and the type of output produced. The option DDFM=SATTERTH requests Satterthwaite's approximation for the denominator degrees of freedom for tests of the fixed effects. Alternatively, the option CHISQ requests that multivariate Wald tests be computed and

compared to the reference chi-squared distribution. The option S (or SOLUTION) requests that the estimates of the fixed effects, and their standard errors, be displayed.

REPEATED <repeated effect> / SUBJECT = effect <options>;

The REPEATED statement is primarily used to distinguish which observations are correlated and which can be regarded as independent of one another. This is achieved with the SUBJECT option, which is used to denote a variable that distinguishes clusters of correlated responses. By including a variable in the SUBJECT option (e.g., a subject identifier), pairs of observations with the same value of that variable are regarded as correlated while pairs of observations with distinct values are regarded as independent.

The REPEATED statement also includes options for specifying assumptions about the nature of the covariance among the errors. This is achieved with the TYPE=<pattern> option (e.g., TYPE=UN specifies an unstructured covariance matrix). A full listing and description of all the possible covariance patterns can be found in the SAS documentation. There are also various options that modify the type of output that is produced. The option R and RCORR print the covariance and correlation matrices, respectively.

Finally, a variable denoted the “repeated effect” can also be included on the REPEATED statement, and this identifies “units within a cluster.” In the context of longitudinal data, the “repeated effect” identifies the measurement occasions. While it is not always necessary to include this variable, failure to do so may have unforeseen consequences when there are vectors of repeated measures of different length and/or when the vector of responses are not in the same order for all subjects. In [Table 5.11](#) the REPEATED statement identifies “time” as the repeated effect. To avoid any potential problems, it is recommended that this variable be included in the REPEATED statement to ensure that the covariance is structured and estimated appropriately.

## **5.10 FURTHER READING**

A useful review of traditional profile analysis, targeted at applied researchers, can be found in Chapter 3 (Section 3.4, pp. 48–52) of Hand and Taylor (1987).

A more detailed discussion of the subtle issues surrounding adjustment for baseline response in the analysis of change can be found in the articles by Lord (1967), Laird (1983), Fitzmaurice (2001), and Glymour et al. (2005), and in Chapter 7 (Section 7.3, pp. 489–496) of Bock (1975).

# Bibliographic Notes

One of the earliest descriptions of traditional profile analysis appeared in an article by Greenhouse and Geisser (1959). Profile analysis is also discussed in detail in Chapter 6 of Johnson and Wichern (2002), Chapters 4 and 5 of Morrison (1990), and Chapters 5 and 6 of Rencher (2002).

## Problems

**5.1** In the National Cooperative Gallstone Study (NCGS), one of the major interests was to study the safety of the drug chenodiol for the treatment of cholesterol gallstones (Schoenfield et al., 1981; Wei and Lachin, 1984). In this study, patients were randomly assigned to high-dose (750 mg per day), low-dose (375 mg per day), or placebo. We focus on a subset of data on patients who had floating gallstones and who were assigned to the high-dose and placebo groups.

In the NCGS it was suggested that chenodiol would dissolve gallstones but, in doing so, might increase levels of serum cholesterol. As a result serum cholesterol (mg/dL) was measured at baseline and at 6, 12, 20, and 24 months of follow-up. Many cholesterol measurements are missing because of missed visits, laboratory specimens were lost or inadequate, or patient follow-up was terminated.

The NCGS serum cholesterol data are stored in an external file: cholesterol.dat

Each row of the data set contains the following seven variables:

Group ID Y<sub>1</sub> Y<sub>2</sub> Y<sub>3</sub> Y<sub>4</sub> Y<sub>5</sub>

*Note:* The categorical variable Group is coded 1 = High-Dose, 2 = Placebo.

**5.1.1** Read the data from the external file and keep it in a “multivariate” or “wide” format.

**5.1.2** Calculate the sample means, standard deviations, and variances of the serum cholesterol levels at each occasion for each treatment group.

**5.1.3** On a single graph, construct a time plot that displays the mean serum cholesterol versus time (in months) for the two treatment groups. Describe the general characteristics of the time trends for the two groups.

**5.1.4** Next read the data from the external file and put the data in a “univariate” or “long” format, with five “records” per subject.

**5.1.5** Assuming an unstructured covariance matrix, conduct an analysis of response profiles. Determine whether the patterns of change over time differ in the two treatment groups.

**5.1.6** Display the estimated  $5 \times 5$  covariance and correlation matrices for the five repeated measurements of serum cholesterol.

**5.1.7** With baseline (month 0) and the placebo group (group 2) as the *reference group*, write out the regression model for mean serum cholesterol that corresponds to the analysis of response profiles in Problem 5.1.5.

**5.1.8** Let  $L$  denote a matrix of known weights and  $\beta$  the vector of linear regression parameters from the model assumed in Problem 5.1.7. The null hypothesis that the patterns of change over time do not differ in the two treatment groups can be expressed as  $H_0: L\beta = 0$ . Describe an appropriate weight matrix  $L$  for this null hypothesis.

**5.1.9** Show how the *estimated* regression coefficients from an analysis of response profiles can be used to construct the time-specific means in the two groups. Compare these estimated means with the sample means obtained in Problem 5.1.2.

**5.1.10** With baseline (month 0) and the placebo group (group 2) as the *reference group*, provide an interpretation for each of the estimated regression coefficients in terms of the effect of the treatments on the patterns of change in mean serum cholesterol.

<sup>†</sup> Readers may find the level of detail in this section challenging; this section can be omitted at first reading without loss of continuity. However, the reader who returns to it will find that it yields important insights about baseline adjustment.

# *Chapter 6*

## *Modeling the Mean: Parametric Curves*

### **6.1 INTRODUCTION**

In the previous chapter we described an approach to modeling longitudinal data that effectively imposed no structure on the underlying mean response trend over time. This approach has some appeal when all subjects are measured at the same set of occasions and the number of measurement occasions is relatively small (e.g., not more than 4 or 5). But as the number of occasions increases and/or when the repeated measures are irregularly timed, analyzing response profiles becomes much less appealing. Even in cases where the number of repeated measures is relatively small, there are two obvious drawbacks of the analysis of response profiles that limit its usefulness for the analysis of longitudinal data. The first is that a statistical test of the null hypothesis of no group  $\times$  time interaction is an omnibus or global test and provides only a broad assessment of whether the mean response profiles are the same in the different groups. If the null hypothesis is rejected, this does not indicate the specific ways in which the mean response profiles differ. As a result, additional analyses are invariably required. Second, by completely ignoring the time-ordering of the repeated measurements, the analysis of response profiles fails to recognize that they can be considered as observations of some continuous, underlying response process over time. The mean response over time can very often be described by relatively simple parametric (e.g., linear or quadratic) or semiparametric (e.g., piecewise linear) curves. From a purely substantive point of view, it is unlikely that the pattern of change in the mean response over the duration of a longitudinal study will be so complicated that its description requires as many parameters as there are measurement occasions. The analysis of response profiles uses a saturated model for the mean response over time, and thereby produces a perfect fit to the observed mean response profile. At first glance this might seem like a desirable feature of any analytic approach; namely that it fits the observed mean responses well. (In fact, not just well, but perfectly!) However, in doing so, the method fails to describe the most salient aspects of the changes in the mean response over time in terms of some pattern that can be given a substantive or theoretical interpretation. In summary, in the analysis of response profiles there is no reduction in complexity.

In contrast, the fitting of parametric or semiparametric curves to longitudinal data can be justified on both substantive and statistical grounds. Substantively, in many longitudinal studies the true underlying mean response process is likely to change over time in a relatively smooth, monotonically increasing or decreasing pattern, at least for the duration of the study. As a result simple parametric or semiparametric curves can be used to describe how the mean response changes over time. From a statistical perspective the fitting of parsimonious models for the mean response will result in statistical tests of covariate effects (e.g., treatment  $\times$  time interactions) that have greater power than in an analysis of response profiles. The reason for the greater power is that the tests of covariate effects focus only on a relatively narrow range of alternative hypotheses. In contrast, the test statistics in the analysis of response profiles disperse their power over a much broader, but in many cases less substantively plausible or relevant, range of alternative hypotheses. For example, when trends in the mean response over time are assumed to be linear, and a linear trend actually provides a reasonable approximation to the true underlying shape of the mean response profile, the resulting tests of time trends and covariate effects will have greater power than the global tests in an analysis of response profiles. Note, however, that the tests based on parametric curves will only be more powerful at detecting changes in the mean response that exhibit a linear trend. They will not be more powerful, however, if the underlying shape of the mean response over time is U-shaped, rather than

linear. Finally, simple parametric curves provide a parsimonious description of changes in the mean response over time in terms of a relatively small number of parameters. The results can be communicated easily to investigators and empirical researchers. In the following two sections, we describe two broad approaches for describing patterns of change in the mean response over time: polynomial trends and linear splines.

## 6.2 POLYNOMIAL TRENDS IN TIME

One widely adopted approach for analyzing longitudinal data is to describe the patterns of change in the mean response over time in terms of simple polynomial trends, for example, linear or quadratic trends. In this approach the means are modeled as an explicit function of time. This approach can handle highly unbalanced designs in a relatively seamless way. For example, mistimed measurements are easily incorporated in the model for the mean response.

# LINEAR TRENDS OVER TIME

The simplest possible curve for describing changes in the mean response over time is a straight line. In this model the slope for time has direct interpretation in terms of a constant change in the mean response for a single-unit change in time. Consider the hypothetical two-group study comparing a novel *treatment* and a *control* discussed in Section 5.2. If the mean response changes in an approximately linear fashion over the duration of the study, we can adopt the following linear trend model:

$$(6.1) E(Y_{ij}) = \beta_1 + \beta_2 \text{Time}_{ij} + \beta_3 \text{Group}_i + \beta_4 \text{Time}_{ij} \times \text{Group}_i,$$

where  $\text{Group}_i = 1$  if the  $i^{\text{th}}$  individual was assigned to the novel treatment, and  $\text{Group}_i = 0$  otherwise;  $\text{Time}_{ij}$  denotes the measurement time for the  $j^{\text{th}}$  measurement on the  $i^{\text{th}}$  individual. Note that  $\text{Group}_i$  requires only a single index  $i$ , since individuals do not change treatment groups over the course of the study. Also, by using two indices for  $\text{Time}_{ij}$ , we are implicitly allowing for the fact that there may potentially be mistimed measurements (in the latter case,  $\text{Time}_{ij} \neq \text{Time}_{i'j}$ , where  $i$  and  $i'$  denote two different subjects).

In the linear model given by (6.1), the model for the mean for subjects assigned to the control group is

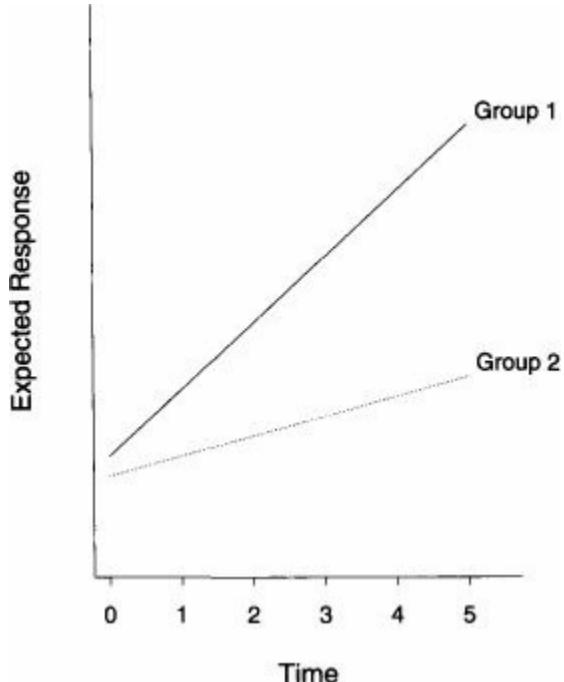
$$E(Y_{ij}) = \beta_1 + \beta_2 \text{Time}_{ij},$$

while for subjects assigned to the treatment group

$$E(Y_{ij}) = (\beta_1 + \beta_3) + (\beta_2 + \beta_4) \text{Time}_{ij}.$$

Thus each group's mean response is assumed to change linearly over time. This model with linear trends for two groups is depicted graphically in [Figure 6.1](#), where the two groups have different intercepts and slopes. Here  $\beta_1$  is the intercept in the control group (the “reference” group), while  $(\beta_1 + \beta_3)$  is the intercept in the treatment group. The intercepts for each of the two groups have interpretation in terms of the mean response when  $\text{Time}_{ij} = 0$ ; more generally,  $\beta_1$  has interpretation as the mean response when all of the covariates are set to zero. Unless some care is taken with how the covariates are scaled (e.g., by centering all quantitative covariates prior to inclusion in the model),  $\beta_1$  is not always readily interpretable and may represent an extrapolation beyond the data at hand. There can also be good reason, beyond issues of parameter interpretation, for centering the variable that denotes the time of measurement; this issue will be discussed later. Finally, the slope, or constant rate of change in the mean response per unit change in time, is  $\beta_2$  in the control group, while the corresponding slope in the treatment group is  $(\beta_2 + \beta_4)$ . Ordinarily, in a longitudinal study the question of primary interest concerns a comparison of the changes in the mean response over time; this can be translated into a comparison of the slopes. Thus, if  $\beta_4 = 0$ , then the two groups do not differ in terms of changes in the mean response over time.

**Fig. 6.1** Graphical representation of model with linear trends for two groups.



The model with linear trend over time is the simplest parametric “curve” that can be used to describe changes in the mean response over time. This model can easily incorporate both discrete (e.g., treatment or exposure group) and quantitative (e.g., dose) covariates. Hypotheses about the dependence of changes in the mean response over time on covariates can be expressed in terms of hypotheses about whether the slope varies as a function of the covariates, that is, in terms of interactions between the covariates and the linear trend in time.

# QUADRATIC TRENDS OVER TIME

When changes in the mean response over time are not linear, higher-order polynomial trends can be considered. For example, if the means are monotonically increasing or decreasing over the course of the study, but in a curvilinear way, a model with quadratic trends can be considered. In a quadratic trend model changes in the mean response are no longer constant (as in the linear trend model) throughout the duration of the study. Instead, the rate of change in the mean response depends on time; that is, the rate of change in the mean response depends on whether the focus is on changes that occur early or later in the study. As a result the rate of change must be expressed in terms of two parameters.

Consider the hypothetical two-group study comparing a novel *treatment* and a *control* discussed earlier. Assuming that the changes in the mean response can be approximated by quadratic trends, the following model can be adopted:

$$\begin{aligned} E(Y_{ij}) &= \beta_1 + \beta_2 \text{Time}_{ij} + \beta_3 \text{Time}_{ij}^2 + \beta_4 \text{Group}_i \\ (6.2) \quad &\quad + \beta_5 \text{Time}_{ij} \times \text{Group}_i + \beta_6 \text{Time}_{ij}^2 \times \text{Group}_i. \end{aligned}$$

In this model the mean response over time for subjects in the control group is given by

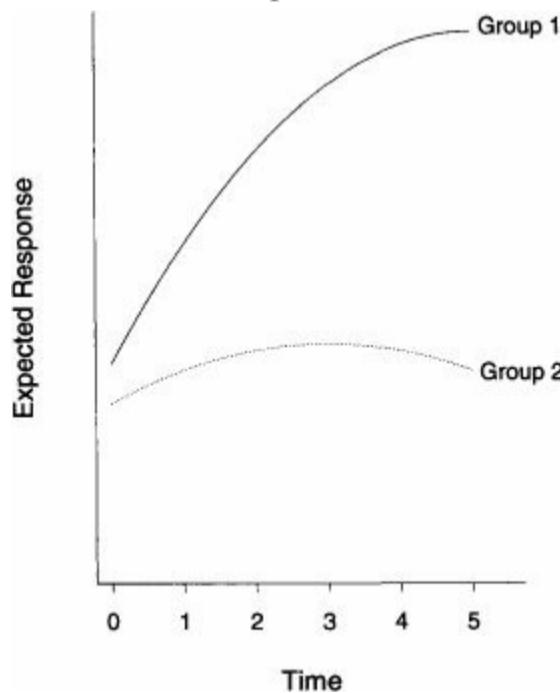
$$E(Y_{ij}) = \beta_1 + \beta_2 \text{Time}_{ij} + \beta_3 \text{Time}_{ij}^2;$$

while the corresponding mean response over time in the novel treatment group is given by

$$E(Y_{ij}) = (\beta_1 + \beta_4) + (\beta_2 + \beta_5) \text{Time}_{ij} + (\beta_3 + \beta_6) \text{Time}_{ij}^2.$$

This model with quadratic trends for two groups is depicted graphically in [Figure 6.2](#), where the two groups have different intercepts (or mean response at time 0) and non-constant rates of change over time that differ between the two groups.

**Fig. 6.2** Graphical representation of model with quadratic trends for two groups.



Note that in the quadratic trends model the mean response changes at a different rate, depending on  $\text{Time}_{ij}$ . For example, the rate of change in the control group is given by  $\beta_2 + 2\beta_3 \text{Time}_{ij}$  (the derivation of this instantaneous rate of change requires some familiarity with calculus and is omitted). Thus early in the study when  $\text{Time}_{ij} = 1$ , the rate of change in the mean response is  $\beta_2 + 2\beta_3$ , whereas later in the study, say  $\text{Time}_{ij} = 4$ , the rate of change in the mean response is  $\beta_2 + 8\beta_3$ . The rate of change is different at the two occasions and the magnitude and sign of the regression coefficients  $\beta_2$  and  $\beta_3$  determine whether the mean response is increasing or decreasing over time and how the rate of change depends on time. The regression coefficients,  $(\beta_2 + \beta_5)$  and  $(\beta_3 + \beta_6)$ , have similar interpretations for the treatment group.

When fitting polynomial trend models, one must take care to avoid extrapolation beyond the data at hand. While polynomial trend models can fit a flexible class of curves to the data, inferences beyond the measurement occasions should be avoided as these will be sensitive to the underlying model

assumptions. While a quadratic trend might be a reasonable approximation for the data, recall that a quadratic trend necessarily has a turning point where the trend changes (e.g., from an increasing trend over time to a decreasing trend, or vice versa). In the absence of a strong theoretical rationale for the model, extrapolation beyond the data can produce nonsensical results and should be avoided.

With polynomial trend models there is a natural hierarchy of effects that has implications for testing hypotheses about linear, quadratic, and higher-order polynomial trends. That is, higher-order terms should be tested (and, if appropriate, removed from the model) before lower-order terms are assessed. Thus in the quadratic model

$$E(Y_{ij}) = \beta_1 + \beta_2 \text{Time}_{ij} + \beta_3 \text{Time}_{ij}^2,$$

it is not meaningful or appropriate to test the coefficient for the linear trend,  $\beta_2$ , in a model that also includes a coefficient for the quadratic trend,  $\beta_3$ . Instead, a test for quadratic trend (versus linear trend) can be performed by testing the null hypothesis that  $\beta_3 = 0$ . If this null hypothesis cannot be rejected, it is then appropriate to remove the quadratic term from the model and consider the model with only linear trend. The test for linear trend is performed by testing the null hypothesis that  $\beta_2 = 0$  in the model that only includes the linear term. This hierarchy is completely analogous to the testing of interactions; that is, tests of main effects (or attempts to interpret the main effects) are not meaningful in the presence of interaction. By the same token, tests of lower-order polynomials in time (e.g., linear trend) are not meaningful in the presence of higher-order polynomials in time (e.g., quadratic trend).

Finally, we return to the issue of “centering” variables. Although “centering”  $\text{Time}_j$  at zero leads to a simple interpretation of the intercept when  $\text{Time}_j$  represents time since baseline, to avoid potential problems of collinearity in the quadratic (or in any higher-order polynomial) trend model, it is advisable to “center”  $\text{Time}_j$  on its mean value. That is, prior to the analysis, replace  $\text{Time}_j$  by its deviation from the mean of  $\text{Time}_1, \text{Time}_2, \dots, \text{Time}_n$ . To highlight the impact of this centering, consider the following example where  $\text{Time}_j \in \{0, 1, 2, \dots, 10\}$ . The correlation between  $\text{Time}_j$  and  $\text{Time}_j^2$  is 0.963. When two covariates are so highly correlated, computational problems associated with collinearity may arise in the estimation of  $\beta$ . Centering  $\text{Time}_j$  before quadratic (or any higher-order polynomial) terms are included in the model helps alleviate problems associated with collinearity. For example, when  $\text{Time}_j \in \{0, 1, 2, \dots, 10\}$  and we create a “centered” variable, say  $\text{CTime}_j = (\text{Time}_j - 5.0)$ , where 5.0 is the mean of  $\{0, 1, 2, \dots, 10\}$ , then the correlation between  $\text{CTime}_j$  and  $\text{CTime}_j^2$  is zero, thereby avoiding any potential problems associated with collinearity. With balanced longitudinal data (requiring only a single index  $j$  for  $\text{Time}_j$ ), it is natural to center  $\text{Time}_j$  on the mean of  $\text{Time}_1, \text{Time}_2, \dots, \text{Time}_n$ . However, with unbalanced longitudinal data, it is important to center  $\text{Time}_{ij}$  at some common value for all individuals. By centering at a common value, the regression intercept is interpretable as the mean response at that common value for time. In general, centering of  $\text{Time}_{ij}$  at individual-specific values (e.g., the mean of the  $n_i$  times of measurement for the  $i^{th}$  individual) should be avoided as these may vary considerably from one individual to another, thereby making the interpretation of the regression intercept meaningless.

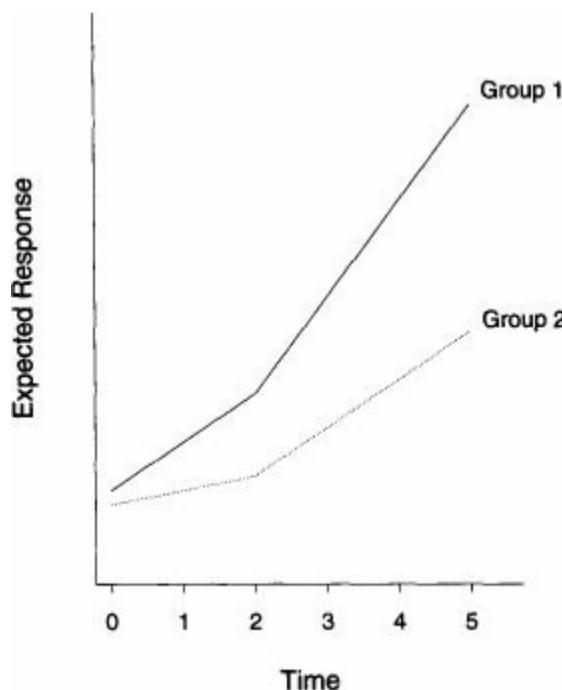
## 6.3 LINEAR SPLINES

Simple parametric curves can provide a parsimonious description of longitudinal trends in the mean response. The simplest case is the linear trend model that characterizes change in the mean response over time in terms of a single slope parameter representing a constant rate of change. By introducing higher-order polynomials in time, various kinds of non-linearities in the longitudinal trends can also be accommodated. However, as the degree of the polynomial increases, the interpretation of the regression coefficients becomes more difficult. As a result the use of polynomials in time is most appealing when any non-linearity can be approximated by quadratic trends.

In some applications the longitudinal trends in the mean response cannot be characterized by first- and second-degree polynomials in time (i.e., linear or quadratic trends). In addition there are other applications where non-linear trends in the mean response cannot be well approximated by polynomials in time of any order. This will most often occur when the mean response increases (or decreases) rapidly for some duration and then more slowly thereafter (or vice versa). When this type of pattern of change arises, it can often be handled by using linear spline models.

If the simplest possible curve is a straight line, then one way to extend the curve is to have a sequence of joined or connected line segments that produces a piecewise linear pattern. Linear spline models provide a very useful and flexible way to accommodate many of the non-linear trends that cannot be approximated by simple polynomials in time. The basic idea behind linear spline models is remarkably simple: divide the time axis into a series of segments and consider a model for the trend over time that is composed of piecewise linear trends, having different slopes within each segment but joined or tied together at fixed times. The locations where the lines meet or are tied together are known as the “knots.” This model allows the mean response to increase or decrease as time proceeds, depending on the sign and magnitude of the regression slopes for the line segments. The resulting piecewise linear curve is called a spline. [Figure 6.3](#) provides an illustration of a linear spline model for two groups with a common knot at time 2. Note that the slopes of the two lines, before and after time 2, are different, with a greater increase in the mean response in the second time segment, and a more attenuated increase in the mean response in the first time segment. This spline model is sometimes referred to as a piecewise linear or “broken-stick” model.

[Fig. 6.3](#) Graphical representation of model with linear splines for two groups, with a common knot at Time = 2.



The simplest possible spline model has only one knot and can be parameterized in a number of different ways. Returning to the hypothetical two-group study comparing a novel *treatment* and a *control* discussed earlier, if the mean response changes over time in a piecewise linear way, we can fit the following linear spline model with knot at  $t^*$ :

$$\begin{aligned} E(Y_{ij}) &= \beta_1 + \beta_2 \text{Time}_{ij} + \beta_3 (\text{Time}_{ij} - t^*)_+ + \beta_4 \text{Group}_i \\ (6.3) \quad &\quad + \beta_5 \text{Time}_{ij} \times \text{Group}_i + \beta_6 (\text{Time}_{ij} - t^*)_+ \times \text{Group}_i, \end{aligned}$$

where  $(x)_+$ , known as a *truncated line function*, is defined as a function that equals  $x$  when  $x$  is positive and is equal to zero otherwise. Thus  $(\text{Time}_{ij} - t^*)_+$  is equal to  $(\text{Time}_{ij} - t^*)$  when  $\text{Time}_{ij} > t^*$  and is equal to zero when  $\text{Time}_{ij} \leq t^*$ . In the model given by (6.3) the means for subjects in the control group are

$$E(Y_{ij}) = \beta_1 + \beta_2 \text{Time}_{ij} + \beta_3 (\text{Time}_{ij} - t^*)_+.$$

When expressed in terms of the mean response prior to and after  $t^*$ ,

$$E(Y_{ij}) = \beta_1 + \beta_2 \text{Time}_{ij}, \quad \text{Time}_{ij} \leq t^*;$$

$$E(Y_{ij}) = (\beta_1 - \beta_3 t^*) + (\beta_2 + \beta_3) \text{Time}_{ij}, \quad \text{Time}_{ij} > t^*.$$

Thus, in the control group, the slope prior to  $t^*$  is  $\beta_2$  and following  $t^*$  is  $(\beta_2 + \beta_3)$ . Similarly the means for subjects in the treatment group are given by

$$E(Y_{ij}) = (\beta_1 + \beta_4) + (\beta_2 + \beta_5) \text{Time}_{ij} + (\beta_3 + \beta_6) (\text{Time}_{ij} - t^*)_+.$$

When expressed in terms of the mean response prior to and after  $t^*$ ,

$$E(Y_{ij}) = (\beta_1 + \beta_4) + (\beta_2 + \beta_5) \text{Time}_{ij}, \quad \text{Time}_{ij} \leq t^*;$$

$$E(Y_{ij}) = \{(\beta_1 + \beta_4) - (\beta_3 + \beta_6) t^*\} \\ + (\beta_2 + \beta_3 + \beta_5 + \beta_6) \text{Time}_{ij}, \quad \text{Time}_{ij} > t^*.$$

Then, in terms of group comparisons, the null hypothesis of no group differences in patterns of change over time can be expressed as  $H_0: \beta_5 = \beta_6 = 0$ . Comparisons of the groups before and after  $t^*$  are also possible. For example, the null hypothesis of no group differences in patterns of change prior to  $t^*$  can be expressed as  $H_0: \beta_5 = 0$ .

The simple spline model considered so far can be extended to include more than one knot or more than two joined line segments. More generally, a spline model with  $K$  knots or break-points will produce  $K + 1$  line segments, and there will be  $K + 1$  corresponding slopes; see Chapter 19 on “smoothing” longitudinal data for a detailed description of splines models with many knots. Thus, in principle, it is possible to accommodate quite complex non-linear patterns for changes in the mean response by including a sufficient number of variables,  $(\text{Time}_{ij} - t_k^*)_+$ , with knots located at  $t_k^*$  (for  $k = 1, \dots, K$ ). However, in practice, the data from many longitudinal studies can be well-approximated by simple piecewise linear models with at most one or two knots that are located at judiciously chosen time points.

Finally, our discussion thus far has avoided the thorny problem of the choice of location(s) for the knot(s). There is an extensive body of research in statistics on automated choices for the knot location, where the location is effectively determined by the data at hand. Ideally the choice of knot location should also incorporate subject-matter considerations. For example, in studies of growth, certain ages are associated with growth spurts. Similarly measures of hormonal response are known to change quite dramatically with the onset of puberty and menopause. In other settings, there may be a body of evidence that the response profile changes in a discernible way at certain time points. An example of the latter is in studies of HIV-infected patients. In early studies of the treatment of HIV-infected patients with AZT, it was increasingly recognized that CD4 counts, a measure of the body’s immune response, increased sharply over a 4- to 6-week period following treatment with AZT, but then leveled off. In summary, the choice of knot location is a mixture of art and science. When it is available, subject-matter knowledge should be brought to bear on the empirical evidence for the most appropriate choice of knot location (see Chapter 19 for a more detailed discussion on choice of knot location).

## 6.4 GENERAL LINEAR MODEL FORMULATION

Below we demonstrate how both polynomial trend and spline models can be expressed in terms of the general linear model

$$E(Y_i|X_i) = \mu_i = X_i\beta,$$

for appropriate choices of  $X_i$ . Let  $n_i$  be the number of repeated measures on the  $i^{th}$  individual ( $i = 1, \dots, N$ ). To illustrate how the polynomial trend model can be expressed in terms of the general linear model, consider the hypothetical two-group study comparing a novel *treatment* and a *control* discussed earlier. Let us assume that the mean response changes over time in a quadratic trend. Then the design matrix  $X_i$  has the following form for the control group:

$$X_i = \begin{pmatrix} 1 & t_{i1} & t_{i1}^2 & 0 & 0 & 0 \\ 1 & t_{i2} & t_{i2}^2 & 0 & 0 & 0 \\ \vdots & \vdots & \vdots & \vdots & \vdots & \vdots \\ 1 & t_{in_i} & t_{in_i}^2 & 0 & 0 & 0 \end{pmatrix},$$

while for the treatment group the design matrix is given by

$$X_i = \begin{pmatrix} 1 & t_{i1} & t_{i1}^2 & 1 & t_{i1} & t_{i1}^2 \\ 1 & t_{i2} & t_{i2}^2 & 1 & t_{i2} & t_{i2}^2 \\ \vdots & \vdots & \vdots & \vdots & \vdots & \vdots \\ 1 & t_{in_i} & t_{in_i}^2 & 1 & t_{in_i} & t_{in_i}^2 \end{pmatrix},$$

where  $t_{ij}$  denotes the time of the  $j^{th}$  measurement on the  $i^{th}$  individual. Then, in terms of the general linear model

$$E(Y_i|X_i) = \mu_i = X_i\beta,$$

where  $\beta = (\beta_1, \dots, \beta_6)'$  is a  $6 \times 1$  vector of regression coefficients, the mean responses in the control group are

$$\mu_i = \begin{pmatrix} \mu_{i1} \\ \mu_{i2} \\ \vdots \\ \mu_{in_i} \end{pmatrix} = \begin{pmatrix} \beta_1 + \beta_2 t_{i1} + \beta_3 t_{i1}^2 \\ \beta_1 + \beta_2 t_{i2} + \beta_3 t_{i2}^2 \\ \vdots \\ \beta_1 + \beta_2 t_{in_i} + \beta_3 t_{in_i}^2 \end{pmatrix},$$

while the mean responses in the treatment group are

$$\mu_i = \begin{pmatrix} \mu_{i1} \\ \mu_{i2} \\ \vdots \\ \mu_{in_i} \end{pmatrix} = \begin{pmatrix} (\beta_1 + \beta_4) + (\beta_2 + \beta_5)t_{i1} + (\beta_3 + \beta_6)t_{i1}^2 \\ (\beta_1 + \beta_4) + (\beta_2 + \beta_5)t_{i2} + (\beta_3 + \beta_6)t_{i2}^2 \\ \vdots \\ (\beta_1 + \beta_4) + (\beta_2 + \beta_5)t_{in_i} + (\beta_3 + \beta_6)t_{in_i}^2 \end{pmatrix}.$$

For the spline model, let us assume that the mean response changes over time in a piecewise linear way, with knot at  $t^* = 4$ . Then the design matrix  $X_i$  has the following form for the control group:

$$X_i = \begin{pmatrix} 1 & t_{i1} & (t_{i1} - 4)_+ & 0 & 0 & 0 \\ 1 & t_{i2} & (t_{i2} - 4)_+ & 0 & 0 & 0 \\ \vdots & \vdots & \vdots & \vdots & \vdots & \vdots \\ 1 & t_{in_i} & (t_{in_i} - 4)_+ & 0 & 0 & 0 \end{pmatrix},$$

while for the treatment group the design matrix is given by

$$X_i = \begin{pmatrix} 1 & t_{i1} & (t_{i1} - 4)_+ & 1 & t_{i1} & (t_{i1} - 4)_+ \\ 1 & t_{i2} & (t_{i2} - 4)_+ & 1 & t_{i2} & (t_{i2} - 4)_+ \\ \vdots & \vdots & \vdots & \vdots & \vdots & \vdots \\ 1 & t_{in_i} & (t_{in_i} - 4)_+ & 1 & t_{in_i} & (t_{in_i} - 4)_+ \end{pmatrix}.$$

The spline model can then be expressed in terms of the general linear model

$$E(Y_i|X_i) = \mu_i = X_i\beta,$$

where  $\beta = (\beta_1, \dots, \beta_6)'$  is a  $6 \times 1$  vector of regression coefficients.

Given that both the polynomial trend and spline models can be expressed in terms of the general linear regression model,

$$E(Y_i|X_i) = \mu_i = X_i\beta,$$

restricted maximum likelihood estimation of  $\beta$ , and the construction of confidence intervals and tests of hypotheses, are possible once the covariance of  $Y_i$  has been specified. Unlike the analysis of response profiles, where the covariance of  $Y_i$  is assumed to be unstructured with no constraints on the covariance parameters other than the requirement that they yield a symmetric matrix (and one that is positive-definite), more parsimonious models for the covariance can be adopted. Indeed, the use of parametric curves for the mean response is most appealing in settings where the longitudinal data are inherently unbalanced over time. As a result an unstructured covariance matrix may not be well-defined, let alone estimated, when, in principle, each individual can have a unique sequence of measurement times. However, the discussion of models for the covariance is postponed until Chapter 7; here we simply assume that some appropriate model for the covariance has been adopted. Given models for both the mean and covariance, REML estimates of  $\beta$ , and their standard errors (based on the estimated covariance of  $\hat{\beta}$ ), can be obtained using the method of estimation described in Chapter 4.

## 6.5 CASE STUDIES

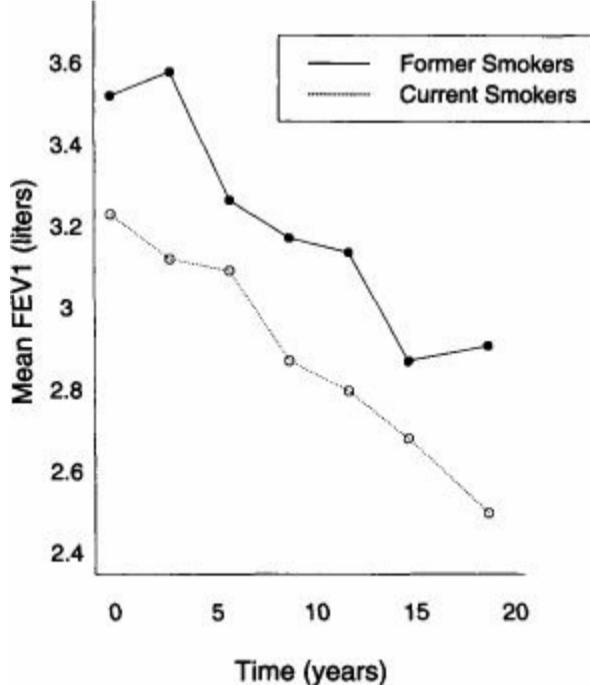
We illustrate the main ideas by considering polynomial trend models for data on lung function ( $FEV_1$ ) from a longitudinal epidemiologic study of current and former smokers aged 36 and older. The application of spline models is illustrated using the blood lead data on the 100 children from the treatment and placebo groups of the Treatment of Lead-Exposed Children (TLC) Trial.

# The Vlagtwedde–Vlaardingen Study

In an epidemiologic study conducted in two different areas in The Netherlands, the rural area of Vlagtwedde in the northeast and the urban, industrial area of Vlaardingen in the southwest, residents were followed over time to obtain information on the prevalence of and risk factors for chronic obstructive lung diseases (van der Lende et al., 1981; Rijcken et al., 1987). Here we focus on a subsample of men and women from the rural area of Vlagtwedde. The sample, initially aged 15 to 44, participated in follow-up surveys approximately every 3 years for up to 21 years. At each survey, information on respiratory symptoms and smoking status was collected by questionnaire and spirometry was performed. Pulmonary function was determined by spirometry and a measure of forced expiratory volume ( $FEV_1$ ) was obtained every three years for the first 15 years of the study, and also at year 19.

In this study,  $FEV_1$  was not recorded for every subject at each of the planned measurement occasions. That is, the data are unbalanced due to incompleteness. The number of repeated measurements of  $FEV_1$  on each subject varied from 1 to 7. For the purpose of this illustration we focus on a subset of the data on 133 residents aged 36 or older at their entry into the study and whose smoking status did not change over the 19 years of follow-up. Each study participant was either a current or former smoker. Current smoking was defined as smoking at least one cigarette per day. The trends in the mean  $FEV_1$  over time, for current and former smokers, are displayed in [Figure 6.4](#). The goal of our analysis is to describe changes in lung function over the 19 years of follow-up with parametric curves and to determine whether the time trends differ for current and former smokers. We summarize differences in mean change between current and former smokers, assuming that the change does not depend strongly on either age or gender (neither variable was available in the data set).

**Fig. 6.4** Mean  $FEV_1$  at baseline (year 0), year 3, year 6, year 9, year 12, year 15, and year 19 in the current and former smoking exposure groups.



First we consider a linear trend in the mean response over time, with intercepts and slopes that differ for the two smoking exposure groups. For all of the analyses reported here, we assume an unstructured covariance matrix. Based on the REML estimates of the regression coefficients in [Table 6.1](#), the mean response for participants who are former smokers is estimated to be

$$E(Y_{ij}) = 3.507 - 0.033 \text{ Time}_{ij},$$

while for participants who are current smokers

$$\begin{aligned} E(Y_{ij}) &= (3.507 - 0.262) - (0.033 + 0.005) \text{ Time}_{ij} \\ &= 3.245 - 0.038 \text{ Time}_{ij}. \end{aligned}$$

Thus it would appear that both groups have a significant decline in mean  $FEV_1$  over time, but there is no discernible difference between the two smoking exposure groups in the constant rate of change,

since the  $\text{Smoke}_i \times \text{Time}_{ij}$  interaction (i.e., the comparison of the two slopes) is not significant, with  $Z = -1.42$ ,  $p > 0.15$ .

**Table 6.1** Estimated regression coefficients and standard errors based on a model with linear trends for the FEV<sub>1</sub> data from the Vlagtwedde–Vlaardingen study.

Variable	Smoking Group	Estimate	SE	Z
Intercept		3.5073	0.1004	34.94
Smoke	Current	-0.2617	0.1151	-2.27
Time		-0.0332	0.0031	-10.84
Smoke × Time	Current	-0.0050	0.0035	-1.42

The adequacy of the linear trend model can be assessed by including higher-order polynomial trends. For example, we can consider a model that allows quadratic trends for changes in FEV<sub>1</sub> over time. Recall that the linear trend model is nested within the quadratic trend model. If the linear trend model is adequate for these data, the difference in maximized log-likelihoods (or the likelihood ratio test statistic) should not be large. The maximized log-likelihoods for the models with linear and quadratic trends are presented in [Table 6.2](#). The likelihood ratio test statistic can be compared to a chi-squared distribution with 2 degrees of freedom (or 6, the number of parameters in the quadratic trend model, minus 4, the number of parameters in the linear trend model). Note that the likelihood ratio test is based on the ML, not REML, log-likelihood because we are comparing two different models for the mean; as a result both models were re-fit using ML estimation. For both models the polynomial trends over time are allowed to differ for the two smoking exposure groups. The likelihood ratio test (comparing the quadratic and linear trend models) produces  $G^2 = 1.3$ , with 2 degrees of freedom ( $p > 0.50$ ). Thus, when compared to the quadratic trend model, the linear trend model appears to be adequate for these data. Finally, for illustrative purposes we can make a comparison with a cubic trend model. This produces a likelihood ratio test statistic,  $G^2 = 4.4$ , with 4 degrees of freedom ( $p > 0.35$ ), indicating again that the linear trend model is adequate for these data.

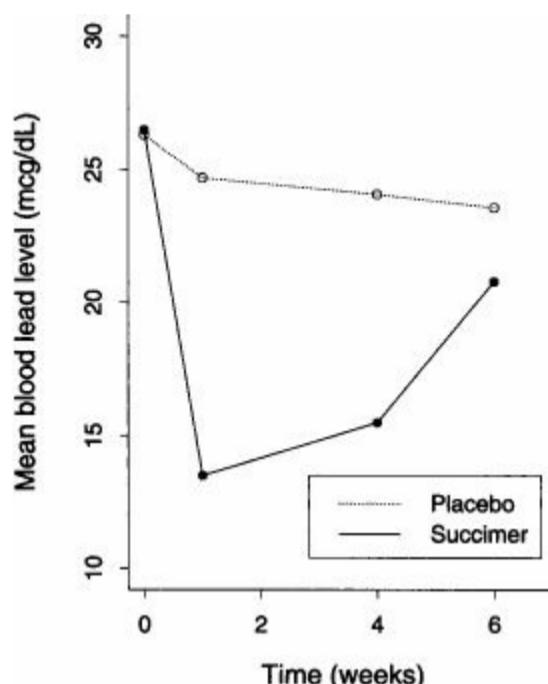
**Table 6.2** Comparison of the maximized (ML) log-likelihoods for the model with linear and quadratic trends for the FEV<sub>1</sub> data from the Vlagtwedde–Vlaardingen study.

Model	-2 (ML) Log-Likelihood
Quadratic Trend Model	237.2
Linear Trend Model	238.5
$-2 \times \text{Log-Likelihood Ratio: } G^2 = 1.3, 2 \text{ df, } (p > 0.50)$	

# Treatment of Lead-Exposed Children Trial

Recall that the TLC trial was a placebo-controlled, randomized trial of a chelating agent, succimer, in children with confirmed blood lead levels of 20 to 44  $\mu\text{g}/\text{dL}$ . The children in the trial were aged 12 to 33 months and lived in deteriorating inner city housing. The following analyses are based on data on blood lead levels at baseline (or week 0), week 1, week 4, and week 6. Note that from the plot of the means in [Figure 6.5](#), it would appear that only the mean blood lead levels in the placebo group can be described by a linear trend; the mean in the succimer group decreases from baseline to week 1, but then increases thereafter.

**Fig. 6.5** Mean blood lead levels at baseline (week 0), week 1, week 4, and week 6 in the succimer and placebo groups.



Given that there are non-linearities in the trends over time, higher-order polynomial models (e.g., a quadratic trend model) could be fit to the data. However, to illustrate the application of spline models, we accommodate the non-linearity with a piecewise linear model with common knot at week 1,

$$E(Y_{ij}) = \beta_1 + \beta_2 \text{Week}_{ij} + \beta_3 (\text{Week}_{ij} - 1)_+ + \beta_4 \text{Group}_i \times \text{Week}_{ij} \\ + \beta_5 \text{Group}_i \times (\text{Week}_{ij} - 1)_+,$$

where  $\text{Group}_i = 1$  if assigned to succimer, and  $\text{Group}_i = 0$  otherwise. Because of the randomization of children to the two treatment groups, the model does not contain a main effect of Group, and we assume a common mean blood lead level at baseline. In this piecewise linear model, the means for subjects in the placebo group are given by

$$E(Y_{ij}) = \beta_1 + \beta_2 \text{Week}_{ij} + \beta_3 (\text{Week}_{ij} - 1)_+,$$

while in the succimer group the means are given by

$$E(Y_{ij}) = \beta_1 + (\beta_2 + \beta_4) \text{Week}_{ij} + (\beta_3 + \beta_5) (\text{Week}_{ij} - 1)_+.$$

The REML estimates of the regression parameters from the piecewise linear model are given in [Table 6.3](#). When expressed in terms of the mean response prior to and after week 1, the estimated means in the placebo group are

**Table 6.3** Estimated regression coefficients and standard errors based on a piecewise linear model, with common knot at week 1, for the blood lead level data from the TLC trial.

Variable	Group	Estimate	SE	Z
Intercept		26.3422	0.4991	52.78
Week		-1.6296	0.7818	-2.08
(Week - 1) <sub>+</sub>		1.4305	0.8777	1.63
Group × Week	S	-11.2500	1.0924	-10.30
Group × (Week - 1) <sub>+</sub>	S	12.5822	1.2278	10.25

$$\hat{\mu}_{ij} = \hat{\beta}_1 + \hat{\beta}_2 \text{Week}_{ij}, \quad \text{Week}_{ij} \leq 1;$$

$$\hat{\mu}_{ij} = (\hat{\beta}_1 - \hat{\beta}_3) + (\hat{\beta}_2 + \hat{\beta}_3) \text{Week}_{ij}, \quad \text{Week}_{ij} > 1.$$

Thus in the placebo group the slope prior to week 1 is  $\hat{\beta}_2 = -1.63$  and, following week 1, is  $(\hat{\beta}_2 + \hat{\beta}_3) = -1.63 + 1.43 = -0.20$ . Similarly, when expressed in terms of the mean response prior to and after week 1, the estimated means for subjects in the succimer group are given by

$$\hat{\mu}_{ij} = \hat{\beta}_1 + (\hat{\beta}_2 + \hat{\beta}_4) \text{Week}_{ij}, \quad \text{Week}_{ij} \leq 1;$$

$$\begin{aligned} \hat{\mu}_{ij} = & \hat{\beta}_1 - (\hat{\beta}_3 + \hat{\beta}_5) \\ & + (\hat{\beta}_2 + \hat{\beta}_3 + \hat{\beta}_4 + \hat{\beta}_5) \text{Week}_{ij}, \quad \text{Week}_{ij} > 1. \end{aligned}$$

The estimates of the mean blood lead levels for the placebo and succimer groups are presented in [Table 6.4](#). The estimated means from the piecewise linear model appear to adequately fit the observed mean response profiles for the two treatment groups.

**Table 6.4** Estimated mean blood lead levels for the placebo and succimer groups from the piecewise linear model, with common knot at week 1; the observed means are in parentheses.

Group	Week 0	Week 1	Week 4	Week 6
Succimer	26.3 (26.5)	13.5 (13.5)	16.9 (15.5)	19.1 (20.8)
Placebo	26.3 (26.3)	24.7 (24.7)	24.1 (24.1)	23.7 (23.6)

Note that the model with linear trends (and common intercept) is nested within the piecewise linear model (since the former can be obtained by setting  $\beta_3 = \beta_5 = 0$  in the latter). When these two models are compared in terms of their maximized log-likelihoods (see [Table 6.5](#)), the likelihood ratio test statistic is  $G^2 = 121.8$  and can be compared to a chi-squared distribution with 2 degrees of freedom (or 5, the number of parameters in the linear spline model, minus 3, the number of parameters in the linear trend model). Note that the likelihood ratio test is based on the ML, not REML, log-likelihood because we are comparing two different models for the mean; as a result both models were re-fit using ML. The magnitude of the likelihood ratio test statistic (with  $p < 0.0001$ ) indicates that the piecewise linear model significantly improves the overall fit to the mean response over time when compared to a linear trend model. This simply confirms what was already obvious from the plot of the means in [Figure 6.5](#). Although the piecewise linear and quadratic trend models (with common intercept for the two treatment groups) are not nested, they both have the same number of parameters and therefore their respective log-likelihoods can be directly compared (see [Table 6.5](#)). From a comparison of the maximized log-likelihoods it is apparent that the piecewise linear model fits these data discernibly better than the quadratic trend model ( $-2 \text{ ML log-likelihood} = 2436.2$  for the piecewise linear model versus  $-2 \text{ ML log-likelihood} = 2511.7$  for the quadratic trend model).

**Table 6.5** Comparison of the maximized (ML) log-likelihoods for the models with linear and quadratic trends, and piecewise linear trend with common knot at week 1, for the blood lead level data from the TLC trial.

Model	$-2 \text{ (ML) Log-Likelihood}$
Piecewise Linear (Spline) Model	2436.2
Quadratic Trend Model	2511.7
Linear Trend Model	2558.0

Finally, it is quite instructive to examine the estimated unstructured covariance matrix for the linear trend model, a model that does not fit these data well. In [Table 6.6](#) the REML estimates of the unstructured covariance matrix are presented. Note that the estimated variances at weeks 1 and 4 are approximately three to four times greater than at baseline. Moreover the estimated covariance matrix

is discernibly different from that obtained in the analysis of response profiles in Section 5.4 that placed no structure on the means. The inflation of the variance at weeks 1 and 4 is indirectly an indication that the lack of fit to the means at these two occasions is being attributed to error variability, and hence the inflation of the variance at these two occasions. This highlights an important issue that will be discussed in greater detail in Chapter 7, namely that there is an interdependence between the mean response and the covariance that has important implications for how these two aspects of longitudinal data are jointly modeled.

**Table 6.6** Estimated unstructured covariance matrix for the linear trend model for the blood lead levels at baseline, week 1, week 4, and week 6 for the children from the TLC trial.

Covariance Matrix			
25.5	13.8	16.1	21.4
13.8	111.2	81.2	38.4
16.1	81.2	78.3	36.8
21.4	38.4	36.8	59.4

In conclusion, one of the main aspects of summarizing trends in the mean response over time via parametric curves is that trends over time and their relation to covariates can be expressed as a function of a small number of parameters. That is, covariate effects on changes in the mean response over time can be captured in one or two regression parameters, leading to more powerful tests when the models are appropriate for the data at hand. Also the parametric curves define the conditional mean of  $Y_{ij}$ ,  $E(Y_{ij}|X_{ij})$ , as an explicit function of the times of measurement,  $\text{Time}_{ij}$ . As a result there is no reason to require all individuals to have the same set of measurement times, nor even the same number of measurements. In our examples we have used only data sets where subjects are measured at the same set of occasions, but this is because we will require models for the covariance when subjects are measured at arbitrary points in time. Hence examples of this will be given in the next chapter.

## 6.6 COMPUTING: FITTING PARAMETRIC CURVES USING PROC MIXED IN SAS

To fit a linear trend model to data from two or more groups measured repeatedly over time, we can use the illustrative SAS commands given in [Table 6.7](#). Note that this model assumes that the covariance matrix is unstructured. In principle, alternative assumption about the covariance can be considered. Indeed, when the data are unbalanced over time, it will be necessary to consider parametric models for the covariance; this topic will be discussed in greater detail in Chapter 7.

**Table 6.7** Illustrative commands for a linear trend model using PROC MIXED in SAS.

---

```
PROC MIXED;
  CLASS id group t;
  MODEL y=group time group*time/S CHISQ;
  REPEATED t/TYPE=UN SUBJECT=id R RCORR;
```

---

Note that the CLASS statement includes a variable *t*. This variable is an additional copy of the variable *time*. The difference is that while *t* is declared as a categorical variable on the CLASS statement, *time* is not and is treated as a quantitative covariate in the MODEL statement. The reason for having two versions, *time* and *t*, one quantitative and the other categorical, is that it is good practice to include, wherever possible, a REPEATED effect. This ensures that the covariance is estimated correctly when the design is balanced but incomplete due to missingness or when the study is balanced and complete but the repeated measures are not in the same order for each subject in the data set (e.g., this might arise when the data set has previously been sorted on another variable). With unbalanced data it will very often not be possible to include a REPEATED effect; instead, the covariance model will need to be defined explicitly in terms of the times of measurement. A further discussion of this point is postponed until Chapter 7.

Next we present illustrative commands for fitting a quadratic trend model in [Table 6.8](#). The MODEL statement now includes both *time* and *timesqr*, the latter is simply an additional variable that is the square of *time* (i.e.,  $time^2$ ). Note that the MODEL statement includes both main effects of *time* and *timesqr*, and their interactions with *group*.

**Table 6.8** Illustrative commands for a quadratic trend model using PROC MIXED in SAS.

---

```
PROC MIXED;
  CLASS id group t;
  MODEL y=group time timesqr group*time group*timesqr / S CHISQ;
  REPEATED t/TYPE=UN SUBJECT=id R RCORR;
```

---

Finally, we present illustrative commands for fitting spline models. In [Table 6.9](#) we present commands in SAS for fitting a model with a single knot at *time* = 4. The MODEL statement includes *time* and *time\_4*, where *time\_4* is a derived variable for  $(time - 4)_+$ . The latter variable can easily be computed in SAS as

**Table 6.9** Illustrative commands for a spline model, with knot at time = 4, using PROC MIXED in SAS.

---

```
PROC MIXED;
  CLASS id group t;
  MODEL y=group time time_4 group*time group*time_4 / S CHISQ;
  REPEATED t/TYPE=UN SUBJECT=id R RCORR;
```

---

```
time_4 = MAX(time - 4, 0);
```

## **6.7 FURTHER READING**

A concise and clear discussion of how to describe patterns of change over time using polynomial trends can be found in Section 3.5 of the book by Hand and Taylor (1987). A general discussion of splines and piecewise linear regression can be found in Chapter 11 (Section 11.5) of Neter et al. (1996).

# Bibliographic Notes

The use of simple parametric curves to describe changes in the mean response over time has its origins in growth curve analysis. Methods for estimation and testing of growth curves were developed by Wishart (1938), Box (1950), and Rao (1958). Potthoff and Roy (1964) proposed an extension of the repeated measures analysis by MANOVA for growth curves; alternative formulations were developed by Rao (1965), Khatri (1966), and Grizzle and Allen (1969).

An excellent discussion of spline models can be found in Chapter 3 of Ruppert et al. (2003), and the references therein.

## Problems

**6.1** In a study of weight gain (Box, 1950) investigators randomly assigned 30 rats to three treatment groups: treatment 1 was a control (no additive); treatments 2 and 3 consisted of two different additives (thiouracil and thyroxin respectively) to the rats' drinking water. Weight, in grams, was measured at baseline (week 0) and at weeks 1, 2, 3, and 4. Due to an accident at the beginning of the study, data on 3 rats from the thyroxin group are unavailable.

The raw data are stored in an external file: `rat.dat`

Each row of the data set contains the following seven variables:

ID Group Y<sub>1</sub> Y<sub>2</sub> Y<sub>3</sub> Y<sub>4</sub> Y<sub>5</sub>

*Note:* The variable Group is coded 1 = control, 2 = thiouracil, and 3 = thyroxin.

**6.1.1** On a single graph, construct a time plot that displays the mean weight versus time (in weeks) for the three groups. Describe the general characteristics of the time trends for the three groups.

**6.1.2** Read the data from the external file and put the data in a “univariate” or “long” format, with five “records” per subject.

**6.1.3** Assume that the rate of increase in each group is approximately constant throughout the duration of the study. Assuming an unstructured covariance matrix, construct a test of whether the rate of increase differs in the groups.

**6.1.4** On a single graph, construct a time plot that displays the *estimated* mean weight versus time (in weeks) for the three treatment groups from the results generated from Problem 6.1.3.

**6.1.5** Based on the results from Problem 6.1.3, what is the estimated rate of increase in mean weight in the control group (group 1)? What is the estimated rate of increase in mean weight in the thiouracil group (group 2)? What is the estimated rate of increase in mean weight in the thyroxin group (group 3)?

**6.1.6** The study investigators conjectured that there would be an increase in weight, but that the rate of increase would level off toward the end of the study. They also conjectured that this pattern of change may differ in the three treatment groups. Assuming an unstructured covariance matrix, construct a test of this hypothesis.

**6.1.7** Compare and contrast the results from Problems 6.1.3 and 6.1.6. Does a model with only a linear trend in time adequately account for the pattern of change in the three treatments groups? Provide results that support your conclusion.

**6.1.8** Given the results of all the previous analyses, what conclusions can be drawn about the effect of the additives on the patterns of change in weight?

# *Chapter 7*

## *Modeling the Covariance*

### **7.1 INTRODUCTION**

Since one of the defining features of longitudinal data is that they are correlated, we must consider approaches for appropriately modeling the covariance or time dependence among the repeated measures obtained on the same individuals. When an appropriate model for the covariance has been adopted, correct standard errors are obtained and valid inferences about the regression parameters can be made. Accounting for the covariance among repeated measures usually increases efficiency or the precision with which the regression parameters can be estimated; that is, the positive correlation among the repeated measures reduces the variability of the estimate of change over time within individuals. Thus in a longitudinal study the positive correlation among repeated measures can be used to advantage in the study of change over time. In addition, when there are missing data, correct modeling of the covariance is often a requirement for obtaining valid estimates of the regression parameters. In general, failure to take account of the covariance among the repeated measures will result in incorrect estimates of the sampling variability and can lead to misleading scientific inferences.

Longitudinal data present us with two aspects of the data that require modeling: the conditional mean response over time and the conditional covariance among repeated measures on the same individuals. Although these two aspects of the data can, in a certain sense, be modeled separately, they are also interrelated. That is, the choices of models for the mean response and the covariance are interdependent. This interdependence arises because the vector of residuals (observed responses minus fitted responses) depends on the specification of the model for the conditional mean.

Put more formally, the covariance between any pair of residuals, say  $\{Y_{ij} - \mu_{ij}(\beta)\}$  and  $\{Y_{ik} - \mu_{ik}(\beta)\}$ , depends on the model for the conditional mean (i.e., depends on  $\beta$ ). A model for the covariance must be chosen on the basis of some model for the mean response; it represents an attempt to account for the covariance among the residuals that results from a specific model for the mean. A different choice of model for the mean, or moreover any misspecification of the model for the mean, can potentially result in a different choice of model for the covariance. As a result of this interdependence between the models for the mean and covariance, we will need to develop an overall modeling strategy that takes this interdependence into account.

## 7.2 IMPLICATIONS OF CORRELATION AMONG LONGITUDINAL DATA

Before considering approaches for modeling the covariance or correlation among repeated measures, it is worth stepping back and considering some of the implications of the correlation among longitudinal data. First, it should be kept in mind that longitudinal data are not only correlated, for the most part they are also positively correlated. Moreover the positive correlation among repeated measures can be used to advantage in the study of change over time. That is, we can capitalize on the positive correlation among longitudinal data when the main focus of the analysis is on change in the mean response over time.

Consider a simple longitudinal study design where it is of interest to measure change in a health outcome “before” and “after” receiving some health intervention. With only two repeated measures of the outcome, the statistical analysis of these data will focus on the difference score, say  $Y_{i2} - Y_{i1}$ , for each individual. Note that the variability of the difference score is given by

$$\begin{aligned}\text{Var}(Y_{i2} - Y_{i1}) &= \text{Var}(Y_{i1}) + \text{Var}(Y_{i2}) - 2\text{Cov}(Y_{i1}, Y_{i2}) \\ &= \sigma_1^2 + \sigma_2^2 - 2\rho_{12}\sigma_1\sigma_2 \\ &= \sigma_1^2 + \sigma_2^2 - 2\rho_{12}\sigma_1\sigma_2,\end{aligned}$$

where  $\rho_{12}$  is the correlation among the pair of responses,  $Y_{i1}$  and  $Y_{i2}$ . On the other hand, suppose that an alternative study design is adopted to assess the impact of the health intervention. Rather than using a longitudinal design, a cross-sectional design is adopted where study participants are randomly assigned to two groups, a group that receives the intervention and a control group that does not. Then the variance of the difference between the responses of any two individuals, when one individual is randomly selected from the intervention group and the other from the control group, is given by

$$\begin{aligned}\text{Var}(Y_{i2} - Y_{i1}) &= \text{Var}(Y_{i1}) + \text{Var}(Y_{i2}) \\ &= \sigma_1^2 + \sigma_2^2\end{aligned}$$

(where  $Y_{i1}$  and  $Y_{i2}$  now denote the responses from two different individuals from the control and intervention groups, respectively).

Thus, provided that the correlation among repeated measures is positive, the variability of the within-individual differences is always smaller than the variability of the between-individual differences. If in this simple illustration we further assume that the variance of the response is constant (over time in the longitudinal study design, and across groups in the cross-sectional study design), with  $\sigma_1^2 = \sigma_2^2 = \sigma^2$ , then the variance of the within-individual differences is simply  $2\sigma^2(1 - \rho)$ , while the variability of the between-individual differences is  $2\sigma^2$ . The ratio of these two variances provides an index of the precision of within-individual differences when compared to between-individual differences. Their ratio (within-individual variance / between-individual variance) can be expressed as  $(1 - \rho)$ . Thus, when the correlation is relatively large and positive, the variability of the within-individual differences (or within-individual changes) can be substantially smaller than that for the corresponding between-individual differences. It is in this sense that a longitudinal study can provide a more precise (i.e., less variable) estimate of change in the mean response than a cross-sectional study with the same number and pattern of observations.

Finally, it must be emphasized that failure to adequately account for the correlation among repeated measures can result in misleading inferences. For instance, if it is assumed that the repeated measures are uncorrelated, when in fact there is strong positive correlation, the nominal standard errors (resulting from the naive assumption of independent or uncorrelated repeated measures) will be incorrect. Specifically, for contrasts that estimate change in the mean response over time, the nominal standard errors will be too large. In this case, one fails to get the full benefit of longitudinal data. With incorrect standard errors, test statistics and  $p$ -values will also be incorrect and thus can lead to incorrect inferences about patterns of change and their relation to covariates. In addition, when there is missingness that is MAR, but not MCAR, likelihood-based estimation of the regression

parameters,  $\beta$ , requires that the entire joint distribution of the vector of responses be correctly specified. As a result the model for the covariance must be correctly specified to ensure that valid estimates of  $\beta$  are obtained. In general, when there are missing data, greater care must be exercised when modeling the covariance among the responses (see Chapters 17 and 18).

In summary, the positive correlation among repeated measures is an inescapable feature of longitudinal data that must be accounted for in the analysis in order to make appropriate inferences. Although the correlation, or more generally, the covariance among the repeated responses, is not usually of intrinsic interest, it cannot simply be ignored. Moreover the positive correlation among longitudinal data enables us to estimate changes in the mean response, and their relation to covariates, with far greater precision than would be possible if the data were uncorrelated. Recognizing that the covariance is an important aspect of the data that must be properly accounted for to complete the specification of any regression model for longitudinal data, there are three broad approaches to modeling the covariance that can be distinguished. The first is to allow any arbitrary pattern of covariance among the repeated measures; this approach results in an “unstructured” covariance and is the topic of Section 7.3. Alternatively, structure can be placed on the covariance matrix and there are two main strategies for doing so. The first modeling approach borrows ideas from the time series literature and assumes that the variances and covariances are not arbitrary but follow distinctive patterns. As a result we refer to these models as covariance pattern models, and they are the topic of Section 7.4. Finally, in a somewhat less direct way, structure can be imposed on the covariance through the introduction of *random effects* in the model for the mean response. That is, by assuming that the mean response depends on a combination of population parameters,  $\beta$  (also known as fixed effects), and individual-specific random effects, a very distinctive structure can be imposed on the covariance matrix. Because of the important role of the random effects structure in modeling the covariance in longitudinal data, discussion of these models will be the topic of Chapter 8.

## 7.3 UNSTRUCTURED COVARIANCE

When the number of measurement occasions is relatively small and all individuals are measured at the same set of occasions, it may be reasonable to allow the covariance matrix to be arbitrary, with all of its elements unconstrained. The only formal requirement is that the covariance matrix be symmetric and positive-definite (recall that the latter condition ensures that while the repeated measures can be highly correlated, there must be no redundancy; that is, none of the repeated measures can be expressed as a linear combination of the others). When no explicit structure is assumed for the covariance among the repeated measures (other than the homogeneity of covariance across different individuals,  $\text{Cov}(Y_i) = \sum_i = \Sigma$ ), the resulting covariance is referred to as an “unstructured” covariance. The chief advantage of an “unstructured” covariance is that no assumptions are made about the variances and covariances. The absence of restrictions on the variances is especially important since our practical experience with many longitudinal studies has led to the empirical observation that the variances are rarely constant over time. For example, the variability of baseline measurements is often discernibly different from the variability of post-baseline measurements.

With  $n$  measurement occasions, the “unstructured” covariance matrix has  $\frac{n \times (n+1)}{2}$  parameters: the  $n$  variances at each occasion and the  $n \times (n - 1)/2$  pairwise covariances (or correlations),

$$\text{Cov}(Y_i) = \begin{pmatrix} \sigma_1^2 & \sigma_{12} & \dots & \sigma_{1n} \\ \sigma_{21} & \sigma_2^2 & \dots & \sigma_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ \sigma_{n1} & \sigma_{n2} & \dots & \sigma_n^2 \end{pmatrix}.$$

Herein lies one of the potential drawbacks of assuming an unstructured covariance: the number of covariance parameters to be estimated grows rapidly with the number of measurement occasions. For example, when there are three occasions ( $n = 3$ ), the number of covariance parameters is 6 (3 variances and 3 pairwise covariances). However, when  $n = 5$ , the number of covariance parameters has grown to 15, while when  $n = 10$ , the number of covariance parameters is 55 (and may be fast approaching the number of subjects enrolled in some longitudinal studies!). When the number of covariance parameters that need to be estimated is large, relative to the sample size, estimation is likely to be very unstable. Thus the use of an unstructured covariance will be appealing only in cases where the number of subjects,  $TV$ , is large relative to the number of covariance parameters,  $\frac{n \times (n+1)}{2}$ .

Setting aside the issue of the potentially large number of covariance parameters that may need to be estimated, the use of an unstructured covariance matrix is problematic when there are mistimed measurements or, more generally, measurements made at grossly irregular intervals. Even the most carefully designed longitudinal study will frequently suffer from deviations from the measurement protocol, resulting in measurements made at arbitrary, irregularly timed intervals. When this problem arises, as it frequently does in studies in the health sciences, the resulting mistimed repeated measurements cannot be accommodated in an unstructured covariance. Thus, when the longitudinal data are inherently unbalanced and/or when the sample size is not sufficiently large to estimate an unstructured covariance, it is usually desirable to impose some structure on the covariance matrix.

## 7.4 COVARIANCE PATTERN MODELS

When attempting to impose some structure on the covariance, a subtle balance needs to be struck. If too little structure is imposed, there may be too many parameters to estimate with the limited amount of data at hand. This was one of the main drawbacks of the unstructured covariance considered in the previous section; by imposing no structure on the covariance, the number of parameters to be estimated grows rapidly with the number of measurement occasions. In a certain sense any given data set contains but a fixed amount of longitudinal information. If too little structure is imposed on the covariance, there will be too many covariance parameters to be estimated from the limited amount of data available, and this will adversely affect the precision with which the main parameters of interest,  $\beta$ , can be estimated. As a result, imposing too little structure on the covariance can result in weaker inferences concerning  $\beta$ . When structure is imposed on the covariance, it is possible to improve the precision with which  $\beta$  can be estimated. However, if too much structure is imposed, there is a potential risk of model misspecification that could ultimately result in misleading inferences concerning  $\beta$ . Once again, this is the classic trade-off between bias and precision. In modeling the covariance, a balance must be struck between these two competing forces.

Structure can be built into the covariance by adopting a covariance pattern model. Covariance pattern models for longitudinal data have their basis in models for serial correlation that were originally developed for time series data. While time series data have a structure that is somewhat different than longitudinal data, being composed of a small number of replications or individuals (in some cases only a single replication) and a large number of repeated measures, they share a common characteristic: the repeated measures are positively correlated and measures taken closer together in time are expected to be more highly correlated than measures further apart in time. Because there are few, if any, replications, much of the statistical literature on the analysis of time series data has focused on parametric models that can describe the covariance structure among the repeated measures with only a few parameters. Many of the models for time series data result in relatively parsimonious models for the covariance that can also be adopted for longitudinal data. Here we describe some of the most widely used covariance pattern models for longitudinal data. Many of these covariance pattern models are available as options in standard statistical software packages for analyzing longitudinal data (e.g., PROC MIXED in SAS).

# Compound Symmetry

Historically one of the first covariance pattern models used for the analysis of repeated measures data was compound symmetry. With a compound symmetry covariance it is assumed that the variance is constant across occasions, say  $\sigma^2$ , and  $\text{Corr}(Y_{ij}, Y_{ik}) = \rho$  for all  $j$  and  $k$ . That is,

$$\text{Cov}(Y_i) = \sigma^2 \begin{pmatrix} 1 & \rho & \rho & \dots & \rho \\ \rho & 1 & \rho & \dots & \rho \\ \rho & \rho & 1 & \dots & \rho \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ \rho & \rho & \rho & \dots & 1 \end{pmatrix},$$

with the constraint that  $\rho \geq 0$ .

The compound symmetry covariance has a randomization justification in certain repeated measures designs (e.g., split-plot designs). In an experiment where the within-subject factor is randomly allocated to subjects, randomization arguments can be made to show that the constant variance and constant correlation conditions hold. (See Chapter 21 for a more detailed discussion of the randomization argument.) However, the randomization argument is simply not justifiable in the longitudinal data setting since measurement occasions cannot be randomly allocated to subjects.

As mentioned in Chapter 3, the compound symmetry covariance does have a theoretical justification when the mean response is thought to depend on a combination of population parameters,  $\beta$ , and a single individual-specific random effect. When the model for the longitudinal responses is expressed as

$$Y_{ij} = X'_{ij}\beta + b_i + \epsilon_{ij},$$

where  $b_i$  is a random effect and  $\epsilon_{ij}$  is a within-individual measurement error, this induces marginally (or averaged over the random effect) a compound symmetry structure on the covariance matrix (with the constraint that  $\rho \geq 0$ ). A more detailed discussion of random effects structures for the covariance will be given in Chapter 8.

The compound symmetry covariance is very parsimonious, with only two parameters regardless of the number of measurement occasions. However, it does make the rather strong assumption that the correlation between any pair of measurements is the same regardless of the time interval between the measurements. This latter aspect of the compound symmetry covariance, the constraint on the correlation among repeated measurements, is somewhat unappealing for most longitudinal data, where the correlations are expected to decay with increasing separation in time. Also the assumption of constant variance across time is unrealistic in many settings. As mentioned earlier, our practical experience with many longitudinal studies has led to the empirical observation that the variances are rarely constant over time. In many settings the assumption of constant variance is the one that is not valid with longitudinal data.

## Toeplitz

The Toeplitz covariance pattern makes the assumption that any pair of responses that are equally separated in time have the same correlation. When the covariance has a Toeplitz form, it is assumed that the variance is constant across occasions, say  $\sigma^2$ , and  $\text{Corr}(Y_{ij}, Y_{ij+k}) = \rho_k$  for all  $j$  and  $k$ . That is,

$$\text{Cov}(Y_i) = \sigma^2 \begin{pmatrix} 1 & \rho_1 & \rho_2 & \dots & \rho_{n-1} \\ \rho_1 & 1 & \rho_1 & \dots & \rho_{n-2} \\ \rho_2 & \rho_1 & 1 & \dots & \rho_{n-3} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ \rho_{n-1} & \rho_{n-2} & \rho_{n-3} & \dots & 1 \end{pmatrix}.$$

Because a Toeplitz covariance assumes that the correlation among responses at adjacent measurement occasions is constant,  $\rho_1$ , this structure is only appropriate when the measurements are made at equal (or approximately equal) intervals of time. Note that the Toeplitz covariance has  $n$  parameters (1 variance parameter, and  $n - 1$  correlation parameters). A special case of the Toeplitz covariance is the (first-order) autoregressive covariance.

# Autoregressive

In the autoregressive model for the covariance it is assumed that the variance is constant across occasions, say  $\sigma^2$ , and  $\text{Corr}(Y_{ij}, Y_{ij+k}) = \rho^k$  for all  $j$  and  $k$ , and  $\rho \geq 0$ . That is,

$$\text{Cov}(Y_i) = \sigma^2 \begin{pmatrix} 1 & \rho & \rho^2 & \dots & \rho^{n-1} \\ \rho & 1 & \rho & \dots & \rho^{n-2} \\ \rho^2 & \rho & 1 & \dots & \rho^{n-3} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ \rho^{n-1} & \rho^{n-2} & \rho^{n-3} & \dots & 1 \end{pmatrix}.$$

The autoregressive covariance is very parsimonious and has only two parameters, regardless of the number of measurement occasions. Because the autoregressive covariance has a Toeplitz form, this structure is only appropriate when the measurements are made at equal (or approximately equal) intervals of time. Note that the correlations decline over time as the separation between pairs of repeated measures increases. However, as mentioned in Section 2.5, in many settings the correlations among repeated measures on the same individual rarely decay that quickly over time.

The autoregressive covariance has a theoretical justification when the errors,  $e_{ij}$ , are thought of as arising from the following first-order autoregressive process:

$$e_{ij} = \rho e_{ij-1} + w_{ij},$$

where  $w_{ij} \sim N(0, \sigma^2 [1 - \rho^2])$  and the process is initiated by an error, say  $e_{i0}$ , where  $e_{i0} \sim N(0, \sigma^2)$ . The autoregressive process is said to be “first-order” because there is only dependence on the previous error; dependence on the two previous errors would yield a “second-order” autoregressive process. Thus the autoregressive covariance can be thought of as resulting from a process where the error term at the  $j^{th}$  occasion is a deterministic function of the error at the previous occasion,  $\rho e_{i,j-1}$  (i.e., the recent past predicts the present), plus an additional (and independent) source of random error,  $w_{ij}$ . For such a process, it can be shown that

$$\text{Var}(e_{ii}) = \sigma^2$$

and

$$\text{Cov}(e_{ij}, e_{ik}) = \sigma^2 \rho^{|j-k|}.$$

Finally, the compound symmetry, Toeplitz, and autoregressive covariances assume that the variances are constant across time. This assumption can easily be relaxed, and it is possible to consider versions of these three covariance pattern models with heterogeneous variances,  $\text{Var}(Y_{ij}) = \sigma_j^2$ . Thus a heterogeneous (variances) autoregressive covariance pattern model is given by

$$\text{Cov}(Y_i) = \begin{pmatrix} \sigma_1^2 & \rho\sigma_1\sigma_2 & \rho^2\sigma_1\sigma_3 & \dots & \rho^{n-1}\sigma_1\sigma_n \\ \rho\sigma_1\sigma_2 & \sigma_2^2 & \rho\sigma_2\sigma_3 & \dots & \rho^{n-2}\sigma_2\sigma_n \\ \rho^2\sigma_1\sigma_3 & \rho\sigma_2\sigma_3 & \sigma_3^2 & \dots & \rho^{n-3}\sigma_3\sigma_n \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ \rho^{n-1}\sigma_1\sigma_n & \rho^{n-2}\sigma_2\sigma_n & \rho^{n-3}\sigma_3\sigma_n & \dots & \sigma_n^2 \end{pmatrix},$$

and has  $n + 1$  parameters ( $n$  variance parameters and 1 correlation parameter).

## Banded

The banded covariance patterns make the assumption that the correlation is zero beyond some specified interval. For example, a banded covariance pattern with a band size of 3 assumes that  $\text{Corr}(Y_{ij}, Y_{ij+k}) = 0$  for  $k \geq 3$ . It is possible to apply a banded pattern to any of the covariance pattern models considered so far. Thus a banded Toeplitz covariance pattern with a band size of 2 is given by

$$\text{Cov}(Y_i) = \sigma^2 \begin{pmatrix} 1 & \rho_1 & 0 & \dots & 0 \\ \rho_1 & 1 & \rho_1 & \dots & 0 \\ 0 & \rho_1 & 1 & \dots & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & 0 & \dots & 1 \end{pmatrix},$$

where  $\rho_2 = \rho_3 = \dots = \rho_{n-1} = 0$ .

Banding makes a very strong assumption about how quickly the correlation decays to zero with increasing separation between the repeated measurements. In our experience with longitudinal studies in the health sciences, it is rare for the correlation to decay to zero, even in studies where there is a lengthy period of follow-up.

# Exponential

When the measurement occasions are not equally spaced over time, the formulation of the autoregressive covariance model can be generalized as follows: Let  $\{t_{i1}, \dots, t_{in}\}$  denote the observation times for the  $i^{th}$  individual, and assume that the variance is constant across all measurement occasions, say  $\sigma^2$ , and

$$\text{Corr}(Y_{ij}, Y_{ik}) = \rho^{|t_{ij} - t_{ik}|},$$

for  $\rho \geq 0$ . That is, the correlation between any pair of repeated measures decreases exponentially with the time separations between them. This structure is referred to as an “exponential” covariance model because it can be re-expressed as

$$\begin{aligned}\text{Cov}(Y_{ij}, Y_{ik}) &= \sigma^2 \rho^{|t_{ij} - t_{ik}|} \\ &= \sigma^2 \exp(-\theta |t_{ij} - t_{ik}|),\end{aligned}$$

where  $\theta = -\log(\rho)$  or  $\rho = \exp(-\theta)$  for  $\theta \geq 0$ . Also note that the exponential covariance model is invariant under linear transformation of the time scale. If we replace  $t_{ij}$  by  $(a + bt_{ij})$  (e.g., if we replace time measured in “weeks” by time measured in “days”), the same form for the covariance matrix holds.

A distinctive feature of the exponential model is that it assumes that the correlation is one if measurements are made repeatedly at the same occasion (or replicate measurements on an individual can be obtained at the same occasion), and that the correlation decreases rapidly to zero as the time separation between measurements increases. This first aspect of the exponential covariance model corresponds to an assumption that the responses are measured without error, an unrealistic assumption in most longitudinal studies in the health sciences. The latter feature, correlations among repeated measurements that decay to zero, is rarely observed in longitudinal studies.

# Hybrid Models

Finally, by combining the autoregressive and the compound symmetry models, it is possible to overcome many of the unappealing aspects of each of these models for longitudinal data. Consider a model for the covariance where

$$\text{Cov}(Y_i) = \Sigma_1 + \Sigma_2,$$

where

$$\Sigma_1 = \sigma_1^2 \begin{pmatrix} 1 & \rho_1 & \rho_1 & \dots & \rho_1 \\ \rho_1 & 1 & \rho_1 & \dots & \rho_1 \\ \rho_1 & \rho_1 & 1 & \dots & \rho_1 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ \rho_1 & \rho_1 & \rho_1 & \dots & 1 \end{pmatrix}$$

and

$$\Sigma_2 = \sigma_2^2 \begin{pmatrix} 1 & \rho_2^{|t_{i1}-t_{i2}|} & \rho_2^{|t_{i1}-t_{i3}|} & \dots & \rho_2^{|t_{i1}-t_{in}|} \\ \rho_2^{|t_{i2}-t_{i1}|} & 1 & \rho_2^{|t_{i2}-t_{i3}|} & \dots & \rho_2^{|t_{i2}-t_{in}|} \\ \rho_2^{|t_{i3}-t_{i1}|} & \rho_2^{|t_{i3}-t_{i2}|} & 1 & \dots & \rho_2^{|t_{i3}-t_{in}|} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ \rho_2^{|t_{in}-t_{i1}|} & \rho_2^{|t_{in}-t_{i2}|} & \rho_2^{|t_{in}-t_{i3}|} & \dots & 1 \end{pmatrix}.$$

In this model

$$\text{Var}(Y_{ij}) = \sigma_1^2 + \sigma_2^2,$$

$$\text{Cov}(Y_{ij}, Y_{ik}) = \rho_1 \sigma_1^2 + \rho_2^{|t_{ij}-t_{ik}|} \sigma_2^2,$$

$$\text{Corr}(Y_{ij}, Y_{ik}) = \frac{\rho_1 \sigma_1^2 + \rho_2^{|t_{ij}-t_{ik}|} \sigma_2^2}{\sigma_1^2 + \sigma_2^2}.$$

This implies that the correlation between replicate measurements on an individual obtained at the same occasion is

$$\frac{\rho_1 \sigma_1^2 + \sigma_2^2}{\sigma_1^2 + \sigma_2^2},$$

which is less than one when  $\rho_1 < 1$ . Furthermore, as the time separation increases, the correlation no longer decays to zero but has a minimum of

$$\frac{\rho_1 \sigma_1^2}{\sigma_1^2 + \sigma_2^2},$$

which is greater than zero provided  $\rho_1 > 0$ . As noted, the compound symmetry model is also a random effects model, so that  $\Sigma_1$  can be written as

$$\Sigma_1 = \begin{pmatrix} \sigma_b^2 + \sigma_\epsilon^2 & \sigma_b^2 & \sigma_b^2 & \dots & \sigma_b^2 \\ \sigma_b^2 & \sigma_b^2 + \sigma_\epsilon^2 & \sigma_b^2 & \dots & \sigma_b^2 \\ \sigma_b^2 & \sigma_b^2 & \sigma_b^2 + \sigma_\epsilon^2 & \dots & \sigma_b^2 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ \sigma_b^2 & \sigma_b^2 & \sigma_b^2 & \dots & \sigma_b^2 + \sigma_\epsilon^2 \end{pmatrix},$$

so that  $\sigma_1^2 = \sigma_b^2 + \sigma_\epsilon^2$ , and  $\rho_1 = \sigma_b^2 / (\sigma_b^2 + \sigma_\epsilon^2)$ . Thus we can think of the total variance,  $\text{Var}(Y_{ij})$ , as the sum of the autoregressive variance,  $\sigma_2^2$ , subject-to-subject variability,  $\sigma_b^2$ , and measurement error variability,  $\sigma_\epsilon^2$ .

## 7.5 CHOICE AMONG COVARIANCE PATTERN MODELS

As mentioned at the beginning of this chapter, the choices of models for the covariance and for the mean are interdependent. As a result it is important to follow a modeling strategy that will result in a sensible choice of models for both aspects of the data. Since model selection criteria for the mean response depend on the correct specification of the model for the covariance (e.g., confidence intervals and tests of hypotheses concerning components of  $\beta$  depend critically on the correct model for the covariance), the first step is to choose a suitable model for the covariance.

It must be recognized that any model for the covariance depends on the assumed model for the mean. A model for the covariance tries to account for the covariance among the residuals, say  $\{Y_{ij} - \mu_{ij}(\beta)\}$  and  $\{Y_{ik} - \mu_{ik}(\beta)\}$ , that result from a specific model for the mean. Therefore the choice of model for the covariance should be based on a “maximal” model for the mean that minimizes any potential misspecification of the model for the mean. Recall that any misspecification of the model for the mean can result in a certain amount of spurious covariance among the residuals, and can induce spurious dependence of the covariance on the covariates.

In longitudinal studies with balanced designs and a very small number of discrete covariates that can be classified as between-subject factors (e.g., treatment assignments, exposure levels, or some characteristic of the subjects), the choice of maximal model is relatively straightforward, since it is possible to choose as the maximal model one that includes the main effects of time (regarded as a within-subject factor) and all other main effects, in addition to their two-way and higher-way interactions. For example, with  $n$  measurement occasions and a single grouping factor with  $G$  levels (e.g., treatment versus control), it is possible to fit a saturated model for the mean response with separate parameters for the  $G \times n$  means. This corresponds to a model with main effects for both the grouping factor and time, in addition to their interaction. This strategy of fitting saturated models for the mean response will be appropriate for longitudinal studies with balanced designs and where the number of qualitatively different levels of the covariates is relatively small. A saturated model for the mean response allows an arbitrary pattern for the mean response profile at every different level of the covariates and thereby minimizes any potential concerns about the impact of misspecification of the model for the mean.

However, in longitudinal studies where there are many covariates (some of which may be quantitative, rather than discrete), the choice of a maximal model is somewhat more difficult. In this case it is not realistic to consider a saturated model for the mean response; instead, a maximal model should be in a certain sense the most elaborate or complex model for the mean response that we would consider from a subject-matter point of view. Such a model may need to distinguish treatment covariates (e.g., treatment groups in experiments) or quasi-treatment covariates (e.g., exposure groups in observational studies) that are the main focus of the study from other covariates that are regarded as potential confounders or effect modifiers. The maximal model will ordinarily include the main effects of the treatment or quasi-treatment covariates and their interactions with time, since the latter effects characterize how changes in the mean response depend on these covariates. The choice of whether to include additional interactions, and so on, must be made on subject-matter grounds. In summary, when there are many potential covariates that can be included in the model for the mean, it is not straightforward to give a simple prescription for choosing the maximal model. The choice of maximal model, it must be recognized, cannot be made through any automatic procedure but must, rather, reflect substantive subject-matter considerations. The maximal model for the mean is a model that excludes certain higher-order interactions among the potential covariates and usually is more complex than any of the sequence of models for the mean response under consideration from a subject-matter point of view. In a sense the reader should envisage a model that, in its degree of complexity, goes a step beyond any model for which empirical researchers in the field would care to provide a specific rationale. Once a maximal model has been chosen, the residual variation and

covariation can then be used to select an appropriate model for the covariance.

Given a maximal model for the mean, a sequence of covariance pattern models can be fit to the data at hand. The choice among models can be made by comparing the maximized likelihoods for each of the covariance pattern models. That is, when any pair of models is nested, a likelihood ratio test statistic can be constructed that compares the “full” and “reduced” models. Recall that two covariance models are said to be nested when the “reduced” model is a special case of the “full” model, so that, when the reduced model holds, the full model must necessarily hold. For example, the compound symmetry model is nested within the Toeplitz model since, if the compound symmetry model holds, then the Toeplitz model must necessarily hold, with  $\rho_1 = \rho_2 = \dots = \rho_{n-1}$ . The likelihood ratio test for two nested covariance models can be constructed by comparing the maximized REML log-likelihoods, say  $\hat{l}_{\text{full}}$  and  $\hat{l}_{\text{red}}$ , for the full and reduced models, respectively. The use of REML, as an alternative to ML, is preferred because it reduces the well-known finite sample bias in the estimation of the covariance. The likelihood ratio test is obtained by taking twice the difference in the respective maximized REML log-likelihoods,

$$G^2 = 2(\hat{l}_{\text{full}} - \hat{l}_{\text{red}}),$$

and comparing the statistic to percentiles from a chi-squared distribution with degrees of freedom equal to the difference between the number of covariance parameters in the full and reduced models.

In general, the likelihood ratio test provides a valid method for comparing nested models for the covariance. However, in certain cases the likelihood ratio test may not be valid, depending on the nature of the null hypothesis that is being tested. In particular, when the likelihood ratio test is testing a null hypothesis that is “on the boundary of the parameter space,” the usual conditions required for classical likelihood theory no longer apply. What is meant by testing a null hypothesis that is “on the boundary of the parameter space”? This rather technical point is best illustrated by considering variances. Recall that variances cannot be negative, they must be positive. As a result variances are considered to be bounded from 0 to  $\infty$ . Thus a likelihood ratio test of the null hypothesis that a variance is zero is testing a null hypothesis that is “on the boundary of the parameter space” for a variance. One consequence is that the usual null distribution for the likelihood ratio test is no longer valid. That is, the null distribution for the likelihood ratio test is no longer a chi-squared distribution with degrees of freedom equal to the difference between the number of parameters in the full and reduced models; instead, the null distribution is a mixture of chi-squared distributions. To illustrate the problem, consider the following model where the mean depends on a combination of population parameters,  $\beta$ , and a single individual-specific random effect:

$$Y_{ij} = X'_{ij}\beta + b_i + \epsilon_{ij},$$

where  $b_i$  is a random effect and  $\epsilon_{ij}$  is a within-individual measurement error, with variances  $\sigma_b^2$  and  $\sigma_e^2$ , respectively. This model induces marginally a compound symmetry covariance subject to the constraint that

$$\rho = \text{Corr}(Y_{ij}, Y_{ik}) = \frac{\sigma_b^2}{\sigma_b^2 + \sigma_e^2}$$

Note that a test of the null hypothesis,  $H_0: \sigma_b^2 = 0$  versus  $H_A: \sigma_b^2 > 0$ , is equivalent to a test of the null hypothesis,  $H_0: \rho = 0$  versus  $H_A: \rho > 0$ , (subject to the constraint that  $\rho$  is non-negative). Under the null hypothesis, the repeated measures are assumed to be uncorrelated; under the alternative hypothesis, they are assumed to be positively correlated. In both instances, the null hypothesis is testing on the boundary of the parameter space (i.e., testing that a variance is zero or testing that a non-negative correlation is zero). As a result the null distribution of the likelihood ratio test is not a chi-squared distribution with 1 degree of freedom. Instead, it is an equally weighted mixture of chi-squared distributions with 0 and 1 degrees of freedom (a chi-squared distribution with 0 degrees of freedom has all of its mass or probability at zero). Some intuition for why it is a mixture of chi-squared distributions with 0 and 1 degrees of freedom can be obtained by considering the fit of the model to the data under the null hypothesis. When  $H_0: \rho = 0$  is true, the fit of the model to the data is equally likely to show some evidence of positive or negative correlation among the responses due to sampling variability. When there is evidence of positive correlation,  $\hat{\rho}$  will be positive, but when

there is evidence of negative correlation,  $\hat{\rho}$  will be zero (since under the alternative hypothesis,  $H_A: \rho > 0$ ,  $\rho$  is constrained to be non-negative). When  $H_0: \rho = 0$  is true, there is a 50:50 chance that  $\rho > 0$  (or  $\rho = 0$ ). As a result  $\rho$  only makes contributions to the likelihood ratio test statistic approximately half of the time, when  $\rho$  is positive. The distribution of the likelihood ratio test statistic can be thought of as chi-squared with 1 degree of freedom half of the time, when  $\rho$  is positive (and chi-squared with 0 degrees of freedom, when  $\rho$  is zero, the other half of the time).

A more detailed discussion of the null distribution of the likelihood ratio test under non-standard conditions is beyond the scope of this book. However, the reader should be aware that the comparison of models for the covariance can sometimes be a non-standard problem. In general, when testing a null hypothesis that is on the boundary of the parameter space, the usual null distribution for the likelihood ratio test is no longer valid. If this problem is simply ignored, and the standard null distribution is naively used, the resulting  $p$ -value for the likelihood ratio test will be overestimated (i.e., a  $p$ -value that is too large will be obtained). Consequently failure to account for this problem can lead to the selection of a model for the covariance that is too parsimonious. That is, there is a danger that the model for the covariance is too simple and ignores some inherent structure in the covariance. Because it is not straightforward to determine the correct null distribution for the likelihood ratio test in these non-standard settings, we recommend the use of  $\alpha = 0.1$ , instead of  $\alpha = 0.05$ , when judging the statistical significance of the likelihood ratio test. Use of the  $\alpha = 0.1$  level is a somewhat *ad hoc* solution but protects against selection of a model for the covariance that is too parsimonious. Alternatively, for cases where the null distribution is a known 50:50 mixture of chi-squared distributions, the critical values given in [Table C.1](#) in Appendix C can be used (see Section 8.5 for additional discussion of this topic).

Often it is of interest to compare non-nested models for the covariance. To compare non-nested models, an alternative approach is the Akaike Information Criterion (AIC). According to the AIC, given a set of competing models for the covariance, one should select the model that minimizes

$$\begin{aligned} \text{AIC} &= -2(\text{maximized log-likelihood}) + 2(\text{number of parameters}) \\ &= -2(\hat{l} - c), \end{aligned}$$

where  $\hat{l}$  is the maximized REML log-likelihood and  $c$  is the number of covariance parameters. Note that AIC can similarly be defined as selecting the model that minimizes

$$\text{AIC} = -\hat{l} + c,$$

or the model that maximizes

$$\text{AIC} = \hat{l} - c.$$

Although AIC can be defined in a number of different ways, the basic underlying idea behind AIC is to strike a balance between two competing objectives: the covariance model must be sufficiently complex to provide a good fit to the data, but at the same time a premium is attached to a parsimonious model. This is achieved by extracting a penalty for the estimation of each additional covariance parameter. With these definitions of AIC, it can be used to compare models with the same fixed effects (i.e., the same model for the mean), but different models for the covariance. Note that expanding the definition of  $c$  to include the number of fixed effects parameters,  $\beta$ , would not alter the selection of the model for the covariance provided the model for the mean is held constant.

We note that AIC is but one of a variety of different “information criteria” that have been proposed. Another criterion is the Bayesian Information Criterion (BIC). According to the BIC, given a set of competing models for the covariance, one should select the model that minimizes

$$\begin{aligned} \text{BIC} &= -2(\text{maximized log-likelihood}) + \log N^*(\text{number of parameters}) \\ &= -2(\hat{l} - \log \sqrt{N^*} c), \end{aligned}$$

where  $N^*$  is the number of subjects. The BIC is sometimes defined where  $N^*$  is the number of “effective subjects,”  $N$  in the case of ML estimation and  $N - p$  in the case of REML estimation (where  $p$  is the dimension of  $\beta$ ). The main idea underlying BIC requires some understanding of the Bayesian approach to model selection where the objective is to choose the model that has the highest posterior probability (or largest Bayes factor). While this is a legitimate model selection criterion, it must be emphasized that BIC only approximates this Bayesian criterion; furthermore the BIC extracts

a very large penalty for the estimation of each additional covariance parameter. In general, we do not recommend the use of BIC for covariance model selection as it entails a high risk of selecting a model that is too simple or parsimonious for the data at hand.

Finally, as mentioned earlier, inferences about  $\beta$  depend on the correct specification of the model for the covariance. Recall that confidence intervals and tests of hypotheses concerning components of  $\beta$  rely on standard errors that are obtained by substituting the REML estimate of  $\Sigma_i$  in the expression for  $\text{Cov}(\hat{\beta})$  (see [Eq. \(4.5\)](#) in Section 4.2). Any misspecification of the model for the covariance has negligible impact on the estimates of the regression coefficients; that is, the regression parameter estimates are unbiased even when the covariance has been misspecified. However, misspecification of the covariance results in incorrect standard errors, and this can lead to potentially misleading inferences concerning  $\beta$  (e.g., due to confidence intervals that are too narrow or wide and  $p$ -values that are too small or large). Fortunately, in many cases, valid standard errors for  $\hat{\beta}$  can be obtained when there is concern about misspecification of the covariance. In particular, valid standard errors for  $\hat{\beta}$  can be based on the “sandwich” estimator of  $\text{Cov}(\hat{\beta})$ ; these standard errors are robust to any misspecification of the covariance. Although the “sandwich” estimator is more widely used in the marginal models for discrete longitudinal data that are the focus of Chapters 12 and 13, we note that the “sandwich” estimator of  $\text{Cov}(\hat{\beta})$  can be applied also in the linear models for longitudinal continuous data described in Part II. The “sandwich” estimator of  $\text{Cov}(\hat{\beta})$  will be discussed in greater detail in Chapter 13.

## 7.6 CASE STUDY

Next we illustrate the main ideas by considering covariance pattern models for data from a trial examining the effectiveness of two different exercise therapy regimens.

# Exercise Therapy Trial

In this study, subjects were assigned to one of two weightlifting programs to increase muscle strength. In the first program, hereafter referred to as treatment 1, the number of repetitions of the exercises was increased as subjects became stronger. In the second program, hereafter referred to as treatment 2, the number of repetitions was held constant but the amount of weight was increased as subjects became stronger. Measurements of muscle strength were taken at baseline and on days 2, 4, 6, 8, 10, and 12. However, to illustrate some of the main differences among the covariance models considered earlier, we focus only on measures of strength obtained at baseline (or day 0) and on days 4, 6, 8, and 12.

Before considering models for the covariance, it is necessary to choose a maximal model for the mean response. Here, with a balanced design on time and only two groups, we chose the maximal model to be the saturated model for the mean, with a total of 10 parameters for the response profiles for the two treatment groups.

First, we consider an unstructured covariance matrix, with all 15 of its elements unconstrained. The estimated covariance and correlation matrices are displayed in [Tables 7.1](#) and [7.2](#), respectively. Note that the variance is larger by the end of the study when compared to the variance at baseline; this is a characteristic pattern observed in many longitudinal studies. Furthermore, from examination of [Table 7.2](#), the correlations decrease as the time separation between the repeated measures increases.

**Table 7.1** Estimated unstructured covariance matrix for the strength data at baseline (day 0), day 4, day 6, day 8, and day 12 from the exercise therapy trial.

Day	0	4	6	8	12
0	9.668	10.175	8.974	9.812	9.407
4	10.175	12.550	11.091	12.580	11.928
6	8.974	11.091	10.642	11.686	11.101
8	9.812	12.580	11.686	13.990	13.121
12	9.407	11.928	11.101	13.121	13.944

**Table 7.2** Estimated unstructured correlation matrix for the strength data at baseline (day 0), day 4, day 6, day 8, and day 12 from the exercise therapy trial.

Day	0	4	6	8	12
0	1.0000	0.9237	0.8847	0.8437	0.8102
4	0.9237	1.0000	0.9597	0.9494	0.9017
6	0.8847	0.9597	1.0000	0.9577	0.9113
8	0.8437	0.9494	0.9577	1.0000	0.9394
12	0.8102	0.9017	0.9113	0.9394	1.0000

Despite the apparent increase in the variance over time, we consider an autoregressive model for the covariance. This model is very parsimonious, with only two parameters, one describing the variance,  $\sigma^2$ , the other the correlation,  $\rho$ . When a first-order autoregressive model is fit to the data, it results in the following estimates of the variance and correlation parameters,  $\hat{\sigma}^2 = 11.87$  and  $\hat{\rho} = 0.94$ . The resulting estimated pairwise correlations among the five repeated measurements are given in [Table 7.3](#). This model was fit primarily for illustrative purposes; the model is not very appropriate for these data as they are unequally spaced over time (i.e., there is a four-day interval between the first two repeated measures and the last two repeated measures, but all other adjacent repeated measurements were taken two days apart). In order to account for the unequal time interval, an exponential model for the covariance was considered, where

$$\text{Corr}(Y_{ij}, Y_{ik}) = \rho^{|t_{ij} - t_{ik}|}$$

**Table 7.3** Estimated autoregressive correlation matrix for the strength data at baseline (day 0), day 4,

day 6, day 8, and day 12 from the exercise therapy trial.

Day	0	4	6	8	12
0	1.0000	0.9402	0.8839	0.8311	0.7813
4	0.9402	1.0000	0.9402	0.8839	0.8311
6	0.8839	0.9402	1.0000	0.9402	0.8839
8	0.8311	0.8839	0.9402	1.0000	0.9402
12	0.7813	0.8311	0.8839	0.9402	1.0000

for  $t_{i1} = 0$ ,  $t_{i2} = 4$ ,  $t_{i3} = 6$ ,  $t_{i4} = 8$ , and  $t_{i5} = 12$  for all subjects. This resulted in the following estimates of the variance and correlation parameters,  $\hat{\sigma}^2 = 11.87$  and  $\hat{\rho} = 0.98$ . The resulting estimated pairwise correlations among the five repeated measurements are given in [Table 7.4](#). Of note, the declines in the estimated correlations in [Tables 7.3](#) and [7.4](#) are too fast when compared to the corresponding declines in [Table 7.2](#).

**Table 7.4** Estimated exponential correlation matrix for the strength data at baseline (day 0), day 4, day 6, day 8, and day 12 from the exercise therapy trial.

Day	0	4	6	8	12
0	1.0000	0.9169	0.8780	0.8408	0.7709
4	0.9169	1.0000	0.9576	0.9169	0.8408
6	0.8780	0.9576	1.0000	0.9576	0.8780
8	0.8408	0.9169	0.9576	1.0000	0.9169
12	0.7709	0.8408	0.8780	0.9169	1.0000

Next we consider the choice among these covariance pattern models. The maximized REML log-likelihood and AIC for each of the covariance pattern models are displayed in [Table 7.5](#). Note that there is a hierarchy among the models. The autoregressive and exponential models are both nested within the unstructured covariance. That is, if either the autoregressive or exponential model holds, then the unstructured covariance must necessarily hold. Comparisons of the autoregressive and exponential models with the unstructured covariance can be made using (REML) likelihood ratio tests. However, the autoregressive and exponential models are not nested models; indeed, both models have the same number of parameters. As a result any comparison between these two models can be made directly in terms of their maximized log-likelihoods, since any penalty extracted by information criteria will be the same in both cases (e.g., with AIC a penalty of 4 is extracted for the estimation of the two covariance parameters). The likelihood ratio test, comparing the autoregressive and unstructured covariance, yields

**Table 7.5** Comparison of the maximized (REML) log-likelihoods and AIC for the covariance pattern models for the strength data from the exercise therapy trial.

Covariance Pattern Model	-2 (REML) Log-Likelihood	AIC
Unstructured	597.3	627.3
Autoregressive	621.1	625.1
Exponential	618.5	622.5

$$G^2 = 621.1 - 597.3 = 23.8,$$

and can be compared to a chi-squared distribution with 13 (or  $15 - 2$ ) degrees of freedom. On the basis of the likelihood ratio test there is evidence that the autoregressive model does not provide an adequate fit to the covariance, when compared to the unstructured covariance ( $p < 0.05$ ). On the other hand, the likelihood ratio test, comparing the exponential and unstructured covariance, yields

$$G^2 = 618.5 - 597.3 = 21.2,$$

and when compared to a chi-squared distribution with 13 degrees of freedom,  $p > 0.05$ . Thus the exponential covariance provides an adequate fit to the data. Also, in terms of AIC, the exponential model minimizes this criterion.

## **7.7 DISCUSSION: STRENGTHS AND WEAKNESSES OF COVARIANCE PATTERN MODELS**

The defining feature of covariance pattern models is that they attempt to account for all the potential sources of variability that have an impact on the covariance among repeated measures on the same individual. That is, they do not distinguish between-subject and within-subject sources of variability. Covariance pattern models characterize the covariance among longitudinal data with a relatively small number of parameters. Many of the models (e.g., autoregressive, Toeplitz, and banded) are only appropriate when the repeated measurements are obtained at equal intervals and cannot handle irregularly timed measurements. Although there is a large selection of models for the correlations, the choice of models for the variances is somewhat limited. Covariance pattern models either make the strong assumption that the variances are constant over time, or relax this assumption entirely and allow the variances to depend arbitrarily on time.

For the most part, covariance pattern models are appropriate for balanced longitudinal designs, and many require that the repeated measurements are obtained at equal intervals. Although these models can handle imbalance due to missing data at any of the fixed occasions, they are not well suited for modeling data from inherently unbalanced longitudinal designs. With inherently unbalanced designs, many of the covariance pattern models are not well defined. In an attempt to overcome the latter limitation, a few covariance pattern models have been developed that allow for irregularly timed measurements (e.g., the exponential covariance pattern model). In these models the correlation is assumed to depend on the time separation between pairs of repeated measurements. However, a potential problem with these models is that they assume the correlation decays rapidly with increasing time separation and that the correlation between two measurements taken at the same occasion is one. As mentioned earlier, in our experience with longitudinal studies in the health sciences, the correlation among repeated measures rarely exhibits either of these two characteristics. Furthermore, although these covariance pattern models allow the correlation to depend on the time separation between repeated measurements, they do not allow the variances to depend on time. As a result they make the strong and often unrealistic assumption that the variance remains constant over time.

In conclusion, covariance pattern models are appropriate for balanced longitudinal designs and many models require that the repeated measurements are obtained at equal intervals. In general, we do not recommend the use of covariance pattern models that make the strong assumption that the variances are constant over time. As mentioned earlier, our practical experience with many longitudinal studies has led to the empirical observation that the variances are rarely constant over time. Because the assumption of constant variance is the one that is not valid in many settings, we recommend that covariance pattern models with heterogeneous variances, allowing the variances to depend arbitrarily on time, should generally be adopted.

# 7.8 COMPUTING: FITTING COVARIANCE PATTERN MODELS USING PROC MIXED IN SAS

In the following we assume that there is a single group factor and the maximal model is the saturated model for the mean. Different patterns can be fit to the covariance matrix among the residuals, denoted R in PROC MIXED in SAS, by using the TYPE= option on the REPEATED statement. [Table 7.6](#) provides a summary of some of the commonly used covariance pattern models; a full description of all of the options can be found in the SAS documentation.

**Table 7.6** Covariance pattern modeling options using PROC MIXED in SAS.

---

TYPE = <pattern> Specifies the covariance pattern

UN	Unstructured
CS	Compound symmetry
AR(1)	First-order autoregressive
TOEP	Toeplitz
UN(n)	Banded unstructured, with n bands
CSH	Heterogeneous compound symmetry
ARH(1)	Heterogeneous first-order autoregressive

---

For example, to fit an autoregressive model for the covariance we can use the illustrative SAS commands given in [Table 7.7](#). The options R and RCORR on the REPEATED statement request that the estimated covariance matrix (R) and the corresponding correlation matrix be displayed as part of the output. By default, the covariance and correlation matrices are displayed for the first subject and will have row and column dimensions corresponding to the number of repeated measures obtained on the first subject. When the vector of responses on the first subject is incomplete, it may be preferable to display the covariance and correlation matrices for a subject with complete responses. The options R=1, 5, 7 and RCORR=1, 5, 7 request that the estimated covariance and correlation matrices be displayed for the first, fifth, and seventh subjects.

**Table 7.7** Illustrative commands for an autoregressive model using PROC MIXED in SAS.

---

```
PROC MIXED;
  CLASS id group time;
  MODEL y=group time group*time/S CHISQ;
  REPEATED time/TYPE=AR(1) SUBJECT=id R RCORR;
```

---

To fit covariance pattern models to inherently unbalanced data requires the use of the “spatial” covariance pattern options in PROC MIXED. These are covariance pattern models developed for spatial data that are defined in terms of “distances” in two-dimensional space. However, these options can also be used where “distance” (or time separation) is defined along the single dimension of time. For example, to fit an exponential covariance pattern model, the following option is used:

TYPE = SP(EXP)(list)

where list is the name of the variable used to construct “distances” or time separation between repeated measurements. [Table 7.8](#) contains illustrative commands for fitting an exponential covariance pattern model. Note that the variable `ctime` is simply an additional copy of `time` that is treated as a continuous covariate for the purpose of constructing the time separation between repeated measurements.

**Table 7.8** Illustrative commands for an exponential model using PROC MIXED in SAS.

---

```
PROC MIXED;
```

```
CLASS id group time;  
MODEL y=group time group*time/S CHISQ;  
REPEATED time / TYPE=SP(EXP)(ctime) SUBJECT=id R RCORR;
```

---

Finally, when the EMPIRICAL option is included on the PROC MIXED statement standard errors for  $\hat{\beta}$  are based on the “sandwich” estimator of  $\text{Cov}(\hat{\beta})$ . As mentioned earlier, these standard errors are robust to any misspecification of the model for the covariance. The “sandwich” estimator of  $\text{Cov}(\hat{\beta})$  will be discussed in Chapter 13.

## **7.9 FURTHER READING**

Additional discussion of covariance pattern models can be found in Chapter 6, Section 6.2, of Brown and Prescott (1999) and in the tutorial by Littell et al. (2000).

# Bibliographic Notes

Jennrich and Schluchter (1986) describe covariance pattern models for longitudinal data. For a more recent and comprehensive overview of this topic, see the review article by Zimmerman and Nunez-Anton (2001), and the references therein. Finally, Pourahmadi (1999) presents a flexible approach for parametric modeling of the covariance structure.

Altham (1984) discusses the advantages, in terms of increased precision of estimation of the parameters of interest, that can result from fitting a parsimonious model to complex data. Altham's (1984) general discussion of this issue has great relevance for the modeling of the covariance in longitudinal data.

The large-sample distribution theory for testing a null hypothesis that is “on the boundary of the parameter space” (e.g., testing that a variance is zero) is discussed in Miller (1977), Self and Liang (1987), Stram and Lee (1994, 1995), Silvapulle and Silvapulle (1995), Silvapulle (1996), and Verbeke and Molenberghs (2003).

## Problems

**7.1** In a study of dental growth, measurements of the distance (mm) from the center of the pituitary gland to the pterygomaxillary fissure were obtained on 11 girls and 16 boys at ages 8, 10, 12, and 14 (Potthoff and Roy, 1964).

The raw data are stored in an external file: `dental.dat`

Each row of the data set contains the following six variables:

ID Gender Y<sub>1</sub> Y<sub>2</sub> Y<sub>3</sub> Y<sub>4</sub>

*Note:* The categorical (character) variable Gender is coded F = Female, M = Male. The third measure (at age 12) on subject ID = 20 is a potential outlier.

**7.1.1** On a single graph, construct a time plot that displays the mean distance (mm) versus age (in years) for boys and girls. Describe the time trends for boys and girls.

**7.1.2** Read the data from the external file and put the data in a “univariate” or “long” format, with four “records” per subject.

**7.1.3** For the “maximal” model, assume a saturated model for the mean response. Fit the following models for the covariance:

- (a) unstructured covariance
- (b) compound symmetry
- (c) heterogeneous compound symmetry
- (d) autoregressive
- (e) heterogeneous autoregressive

Choose a model for the covariance that adequately fits the data.

**7.1.4** Given the choice of model for the covariance from Problem 7.1.3, treat age (or time) as a categorical variable and fit a model that includes the effects of age, gender, and their interactions. Determine whether the pattern of change over time is different for boys and girls.

**7.1.5** Show how the *estimated* regression coefficients from Problem 7.1.4 can be used to estimate the means in the two groups at ages 8 and 14.

**7.1.6** Given the choice of model for the covariance from Problem 7.1.3, treat age as a continuous variable and fit a model that includes the effects of a linear trend in age, gender, and their interaction. Compare and contrast the results with those obtained in Problem 7.1.4.

**7.1.7** On a single graph, construct a time plot that displays the *estimated* mean distance (mm) versus age (in years) for boys and girls from the results generated from Problem 7.1.6.

**7.1.8** Show how the regression coefficients from Problem 7.1.6 can be used to estimate the means in the two groups at ages 8 and 14.

**7.1.9** Does a model with only a linear trend in age adequately account for the pattern of change in

the two groups?

**7.1.10** The third measure (at age 12) on subject ID = 20 is a potential outlier. Repeat the analyses in Problems 7.1.3, 7.1.4, 7.1.6 and 7.1.9 excluding the third measure on subject ID = 20. Do the substantive conclusions change?

**7.1.11** Given the results of all the previous analyses, what conclusions can be drawn about gender differences in patterns of dental growth?

# *Chapter 8*

## *Linear Mixed Effects Models*

### **8.1 INTRODUCTION**

In Chapters 5 and 6 we introduced models for longitudinal data where changes in the mean response, and their relation to covariates, can be expressed as

$$E(Y_i|X_i) = X_i\beta,$$

and where the primary goal is to make inferences about the population regression parameters,  $\beta$ . In Chapter 7 we described how the specification of this regression model for longitudinal data can be completed by making additional assumptions about the structure of  $\text{Cov}(Y_i|X_i) = \text{Cov}(e_i) = \sum_i$ . In this chapter we consider an alternative, but closely related, approach for analyzing longitudinal data using linear mixed effects models. The underlying premise of linear mixed effects models is that some subset of the regression parameters vary randomly from one individual to another, thereby accounting for sources of natural heterogeneity in the population. That is, individuals in the population are assumed to have their own subject-specific mean response trajectories over time and a subset of the regression parameters are now regarded as being random. The distinctive feature of linear mixed effects models is that the mean response is modeled as a combination of population characteristics,  $\beta$ , that are assumed to be shared by all individuals, and subject-specific effects that are unique to a particular individual. The former are referred to as *fixed effects*, while the latter are referred to as *random effects*. The term *mixed* is used in this context to denote that the model contains both fixed and random effects.

Although the linear mixed effects model assumes that the longitudinal responses depend on a combination of population (or fixed effects) parameters,  $\beta$ , and subject-specific effects, it nonetheless leads to a model for the marginal mean response (averaged over the distribution of the random effects) that can be expressed in the familiar form

$$E(Y_i|X_i) = X_i\beta.$$

However, the introduction of random effects induces covariance among the responses and  $\text{Cov}(Y_i|X_i) = \sum_i$  has a distinctive random effects structure. With the inclusion of random effects, the covariances among the repeated measures can be expressed as functions of time. Unlike the covariance pattern models considered in Chapter 7, which do not distinguish the different sources of variability that have an impact on the covariance, linear mixed effects models explicitly distinguish between-subject and within-subject sources of variability. Moreover the induced random effects covariance structure can often be described with relatively few parameters, regardless of the number and timing of the measurement occasions.

Because linear mixed effects models explicitly distinguish between fixed and random effects, they allow the analysis of between-subject and within-subject sources of variation in the longitudinal responses. In addition it is not only possible to estimate parameters that describe how the mean response changes in the population of interest, it is also possible to predict how individual response trajectories change over time. For example, linear mixed effects models can be used to obtain predictions of individual growth trajectories over time. The latter will be of interest when the focus of inference is on the individual rather than the population of individuals. For example, in the physician-patient context, these predictions can be used to identify those patients who do not respond well to their assigned treatment in a clinical trial.

One very appealing aspect of linear mixed effects models is their flexibility in accommodating any degree of imbalance in longitudinal data, coupled with their ability to account for the covariance

among the repeated measures in a relatively parsimonious way. That is, with linear mixed effects models we do not require the same number of observations on each subject nor that the measurements be taken at the same set of measurement occasions. As a result these models are particularly well suited for analyzing inherently unbalanced longitudinal data. While the regression models for the mean response described in Chapter 6 can also handle unbalanced longitudinal data, the class of covariance pattern models suitable for unbalanced data is very limited.

# Example: Random Intercept Model

Recall that in earlier chapters we encountered the simplest possible case of a linear mixed effects model: the linear model with a randomly varying subject effect. In this model each subject is assumed to have an underlying level of response that persists over time. This subject effect is incorporated in the linear mixed effects model by regarding it as random, yielding the following model

$$(8.1) \quad Y_{ij} = X'_{ij}\beta + b_i + \epsilon_{ij},$$

where  $b_i$  is the random subject effect and the  $\epsilon_{ij}$  are regarded as measurement or sampling errors. Let us examine this simple model more closely. In this model, the response for the  $i^{th}$  subject at the  $j^{th}$  occasion is assumed to differ from the population mean,  $X'_{ij}\beta$ , by a subject effect,  $b_i$ , and a within-subject measurement error,  $\epsilon_{ij}$ . Both the subject effect and the measurement error are assumed to be random, with mean zero, and with variances,  $\text{Var}(b_i) = \sigma^2_b$  and  $\text{Var}(\epsilon_{ij}) = \sigma^2_\epsilon$ , respectively. In addition it is assumed that  $b_i$  and the  $\epsilon_{ij}$  are independent of one another. Note that this model describes the mean response trajectory over time for any individual,

$$E(Y_{ij}|b_i) = X'_{ij}\beta + b_i,$$

in addition to the mean response profile in the population,

$$E(Y_{ij}) = X'_{ij}\beta,$$

where the averaging is over all individuals in the population. We refer to the former as the *conditional* mean of  $Y_{ij}$ , given the subject-specific effect, and the latter as the *marginal* mean of  $Y_{ij}$  (averaged over the distribution of the subject-specific effects,  $b_i$ ). There is potential for confusion in our use of this terminology, however, since in both cases the mean response is conditional also upon the covariates,  $X_{ij}$ ; for notational convenience, we have suppressed the dependence on covariates. The alert reader will also have noticed a small change in notation from the previous chapters. The measurement or sampling errors in (8.1) are denoted by  $\epsilon_{ij}$  (epsilon) not  $e_{ij}$ . This change in notation is intentional and reflects differences in interpretations of  $\epsilon_{ij}$  and  $e_{ij}$ . In previous chapters, the error  $e_{ij}$  represents the deviation of  $Y_{ij}$  from the mean response in the population,  $X'_{ij}\beta$ . In this chapter, the *within-subject* error  $\epsilon_{ij}$  represents the deviation of  $Y_{ij}$ , from the subject-specific mean response,  $X'_{ij}\beta + b_i$ . Put another way, the random errors,  $e_{ij}$ , in previous chapters have now been decomposed into two random components,  $e_{ij} = b_i + \epsilon_{ij}$ , a between-subject component and a within-subject component.

Next consider the interpretation of the parameters in the model given by (8.1). The regression parameters  $\beta$  describe patterns of change in the mean response over time (and their relation to covariates) in the population of interest, while  $b_i$  describes how the trend over time for the  $i^{th}$  individual deviates from the population average. That is,  $b_i$  represents an individual's deviation from the population mean intercept, after the effects of the covariates have been accounted for. Thus, when combined with the fixed effects,  $b_i$  describes the mean response trajectory over time for any individual. This interpretation is often obscured by the use of vector and matrix notation, but is apparent if we express the model given by (8.1) as

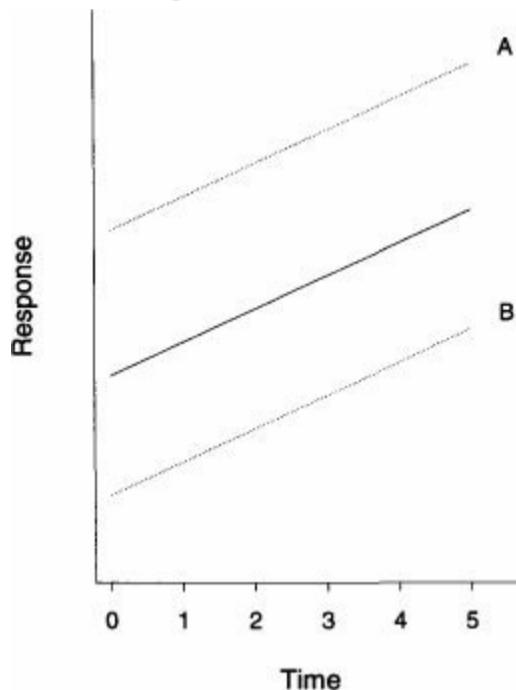
$$\begin{aligned} Y_{ij} &= X'_{ij}\beta + b_i + \epsilon_{ij} \\ &= \beta_1 X_{ij1} + \beta_2 X_{ij2} + \cdots + \beta_p X_{ijp} + b_i + \epsilon_{ij} \\ &= \beta_1 + \beta_2 X_{ij2} + \cdots + \beta_p X_{ijp} + b_i + \epsilon_{ij} \\ &= (\beta_1 + b_i) + \beta_2 X_{ij2} + \cdots + \beta_p X_{ijp} + \epsilon_{ij}, \end{aligned}$$

where  $X_{ij1} = 1$  for all  $i$  and  $j$ , and  $\beta_1$  is then the fixed effect intercept term in the model. When expressed in this way, it can be seen that the intercept for the  $i^{th}$  individual is  $\beta_1 + b_i$  and varies randomly from one individual to another. Because the mean of the random effect  $b_i$  is assumed to be zero,  $b_i$  represents the deviation of the  $i^{th}$  individual's intercept ( $\beta_1 + b_i$ ) from the population intercept,  $\beta_1$ .

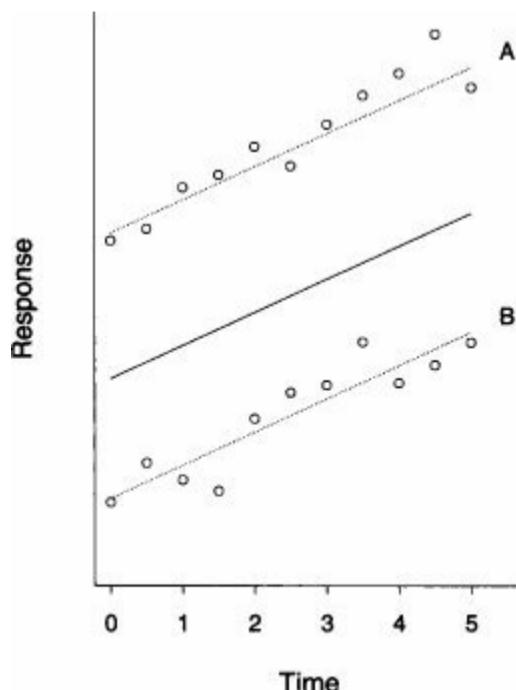
For this simple example of a linear mixed effects model the fundamental ideas can be best

understood by considering the graphical representation of the model equations. [Figure 8.1](#) displays how the marginal mean response over time in the population changes linearly with time (denoted by the solid line), but also indicates how the conditional mean responses for two specific individuals, say subjects A and B, deviate from the population trend (denoted by the broken lines). In this simple illustration, individual A responds “higher” than the population average and thus has a positive  $b_i$ . On the other hand, individual B responds “lower” than the population average and has a negative  $b_i$ . Note that the mixed effects model with randomly varying intercepts does not posit that the repeated measures for individual A or B fall perfectly along these subject-specific response trajectories (represented by the broken lines in [Figure 8.1](#)). The inclusion of the measurement errors,  $\epsilon_{ij}$ , allows the response at any occasion to vary randomly above and below the subject-specific trajectories; this is illustrated in [Figure 8.2](#).

**Fig. 8.1** Graphical representation of the marginal and conditional mean responses over time.



**Fig. 8.2** Graphical representation of the marginal and conditional mean responses over time, plus measurement errors.



Next consider the marginal covariance among the repeated measurements on the same individual. When averaged over the individual-specific effects, the marginal mean of  $Y_{ij}$  is given by

$$E(Y_{ij}) = \mu_{ij} = X'_{ij}\beta.$$

The marginal covariance among the  $Y_{ij}$  is defined in terms of deviations of  $Y_{ij}$  from the marginal mean,  $\mu_{ij}$ . For example, in [Figure 8.2](#) these deviations are positive at all measurement occasions for individual A and negative at all measurement occasions for individual B, indicating a strong positive correlation (marginally) among the responses over time. For the model with randomly varying

intercepts, the marginal variance of each response is given by

$$\begin{aligned}\text{Var}(Y_{ij}) &= \text{Var}(X'_{ij}\beta + b_i + \epsilon_{ij}) \\ &= \text{Var}(b_i + \epsilon_{ij}) \\ &= \text{Var}(b_i) + \text{Var}(\epsilon_{ij}) \\ &= \sigma_b^2 + \sigma^2.\end{aligned}$$

Similarly the marginal covariance between any pair of responses,  $Y_{ij}$  and  $Y_{ik}$ , is given by

$$\begin{aligned}\text{Cov}(Y_{ij}, Y_{ik}) &= \text{Cov}(X'_{ij}\beta + b_i + \epsilon_{ij}, X'_{ik}\beta + b_i + \epsilon_{ik}) \\ &= \text{Cov}(b_i + \epsilon_{ij}, b_i + \epsilon_{ik}) \\ &= \text{Cov}(b_i, b_i) \\ &= \text{Var}(b_i) \\ &= \sigma_b^2.\end{aligned}$$

Thus the marginal covariance matrix of the repeated measurements has the following compound symmetry pattern:

$$\text{Cov}(Y_i) = \begin{pmatrix} \sigma_b^2 + \sigma^2 & \sigma_b^2 & \sigma_b^2 & \dots & \sigma_b^2 \\ \sigma_b^2 & \sigma_b^2 + \sigma^2 & \sigma_b^2 & \dots & \sigma_b^2 \\ \sigma_b^2 & \sigma_b^2 & \sigma_b^2 + \sigma^2 & \dots & \sigma_b^2 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ \sigma_b^2 & \sigma_b^2 & \sigma_b^2 & \dots & \sigma_b^2 + \sigma^2 \end{pmatrix}.$$

This is the only covariance model that arises in both the patterned (see Section 7.4) and random effects families.

Given that the covariance between any pair of repeated measurements is  $\sigma_b^2$ , the correlation is

$$\text{Corr}(Y_{ij}, Y_{ik}) = \frac{\sigma_b^2}{\sigma_b^2 + \sigma^2}.$$

This simple expression for the correlation emphasizes an important aspect of mixed effects models: the introduction of a random subject effect,  $b_i$ , can be seen to induce correlation among the repeated measurements. Although the randomly varying intercepts model is the simplest example of a linear mixed effects model, and the resulting covariance structure is not usually appropriate for longitudinal data, the basic ideas can be generalized to provide a very versatile model for analyzing longitudinal data.

## 8.2 LINEAR MIXED EFFECTS MODELS

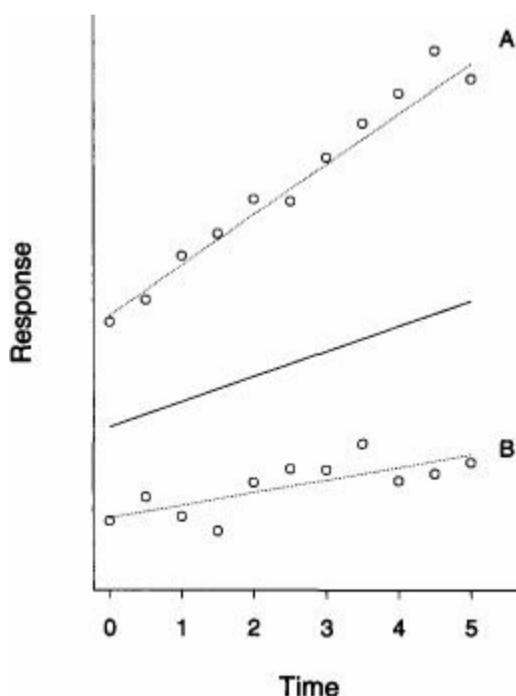
In this section we consider generalizations of (8.1) by allowing additional regression coefficients to vary randomly. We also highlight some of the appealing aspects of the linear mixed effects model alluded to earlier. The underlying premise of the model is that some subset of the regression coefficients vary randomly from one individual to another. In the simplest case considered above, we assumed that the intercept varied randomly. The introduction of this single random effect induces covariance among the repeated measures, albeit with a somewhat restricted form. By allowing a subset of the regression coefficients to vary randomly, a very flexible, and yet quite parsimonious, class of random effects covariance structures becomes available.

To fix ideas, consider the following example of a linear mixed effects model with intercepts and slopes that vary randomly among individuals. That is, for the  $i^{th}$  subject at the  $j^{th}$  measurement occasion,

$$Y_{ij} = \beta_1 + \beta_2 t_{ij} + b_{1i} + b_{2i} t_{ij} + \epsilon_{ij}, \quad j = 1, \dots, n_i.$$

In this model each subject varies not only in their baseline level of response (when  $t_{i1} = 0$ ) but also in terms of changes in their responses over time. This can be best understood by considering the graphical representation of the model equations. [Figure 8.3](#) displays how the marginal mean response in the population changes linearly with time (denoted by the solid line), but also indicates how the conditional mean responses for two specific individuals, say subjects A and B, deviate from the population trend (denoted by the broken lines). In this simple illustration, individual A has a “higher” baseline level of response ( $\beta_1 + b_{1i}$ ) than the population average ( $\beta_1$ ) and thus has a positive  $b_{1i}$ . On the other hand, individual B has a “lower” baseline level of response than the population average and thus has a negative  $b_{1i}$ . In addition individual A has a steeper rate of increase over time ( $\beta_2 + b_{2i}$ ) than the population average ( $\beta_2$ ) and thus has a positive  $b_{2i}$ . Individual B has a less steep rate of increase over time than the population average and thus has a negative  $b_{2i}$ . Finally, the inclusion of the measurement errors,  $\epsilon_{ij}$ , allows the response at any occasion to vary randomly above and below the subject-specific trajectories. In this illustration there are randomly varying intercepts and slopes. However, the linear mixed effects model can be generalized to incorporate additional randomly varying regression coefficients and to allow the means of the random effects to depend on covariates.

[Fig. 8.3](#) Graphical representation of the marginal and conditional mean responses over time, plus measurement errors.



In the following we assume there are  $N$  individuals on whom we have collected  $n_i$  repeated observations, with the response variable  $Y_{ij}$  measured at time  $t_{ij}$ . Thus the longitudinal data can be inherently unbalanced over time. In the most extreme case, each individual has a unique sequence of measurement occasions,  $t_{i1}, \dots, t_{in_i}$ . Although no longitudinal study would ever be intentionally

designed in this way, a change in the metamer for “time” may induce such a design. For example, a longitudinal design can be perfectly balanced when time is defined relative to the baseline measurement but become highly unbalanced if time is defined relative to some landmark event (e.g., puberty, menarche, or menopause).

Using vector and matrix notation, the linear mixed effects model can be expressed as

$$(8.2) \quad Y_i = X_i\beta + Z_i b_i + \epsilon_i,$$

where  $\beta$  is a  $(p \times 1)$  vector of fixed effects,  $b_i$  is a  $(q \times 1)$  vector of random effects,  $X_i$  is a  $(n_i \times p)$  matrix of covariates, and  $Z_i$  is a  $(n_i \times q)$  matrix of covariates, with  $q \leq p$ . Here  $Z_i$  is a known design matrix linking the vector of random effects  $b_i$  to  $Y_i$ . In particular, in many models for longitudinal analysis the columns of  $Z_i$  are a subset of the columns of  $X_i$ . The reason for this restriction on the columns of  $Z_i$  will become evident in Section 8.4; in Chapter 19 we will encounter linear mixed effects models where the columns of  $Z_i$  are no longer a subset of the columns of  $X_i$ . In model (8.2) the particular subset of the regression parameters  $\beta$  that vary randomly is determined by the columns of  $X_i$  that comprise  $Z_i$ . That is, any component of  $\beta$  can be allowed to vary randomly by simply including the corresponding column of  $X_i$  in  $Z_i$ , the design matrix for the random effects. The random effects,  $b_i$ , are assumed to be independent of the covariates,  $X_i$ , and to have a multivariate normal distribution with mean zero and covariance matrix  $G$ . That is,  $E(b_i) = 0$  and  $\text{Cov}(b_i) = G$ . In principle, any multivariate distribution for  $b_i$  could be assumed; in practice,  $b_i$  are assumed to have a multivariate normal distribution.

If, in model (8.2), the vector of random effects,  $b_i$ , has mean zero, the random effects then have interpretation in terms of how the subset of regression parameters for the  $i^{th}$  individual deviate from those in the population. As mentioned previously, the particular subset of the regression parameters,  $\beta$ , that are assumed to vary randomly is determined by the columns of  $X_i$  that comprise  $Z_i$ . For example, in a model with only randomly varying intercepts,  $Z_i$  is a  $(n_i \times 1)$  vector composed of 1's (since  $X_{ij1} = 1$  for all  $i$  and  $j$ ). Later we will consider the form of the design matrix  $Z_i$  for more general models.

An important distinction in the linear mixed effects model is that between the conditional and marginal means of  $Y_{ij}$ . The *conditional* or *subject-specific* mean of  $Y_i$ , given  $b_i$ , is

$$E(Y_i|b_i) = X_i\beta + Z_i b_i,$$

while the *marginal* or population-averaged mean of  $Y_i$ , when averaged over the distribution of the random effects  $b_i$ , is

$$\begin{aligned} E(Y_i) &= \mu_i \\ &= E\{E(Y_i|b_i)\} \\ &= E(X_i\beta + Z_i b_i) \\ &= X_i\beta + Z_i E(b_i) \\ &= X_i\beta, \end{aligned}$$

since  $E(b_i) = 0$ . Thus, in the linear mixed effects model, the vector of regression parameters  $\beta$  (the *fixed effects*) are assumed to be the same for all individuals and have population-averaged interpretations, for example, in terms of changes in the mean response, averaged over all individuals in the population. In contrast to  $\beta$ , the vector  $b_i$  (when combined with the corresponding fixed effects) is composed of subject-specific regression coefficients. These are the *random effects*, and when combined with the fixed effects, they describe the mean response profile of any *individual*. That is, the mean response profile for the  $i^{th}$  individual is given by

$$E(Y_i|b_i) = X_i\beta + Z_i b_i.$$

Finally, the  $(n_i \times 1)$  vector of errors,  $\epsilon_i$ , is assumed to be independent of  $b_i$ , and to also have a multivariate normal distribution with mean zero and covariance matrix  $R_i$ . Ordinarily it is further assumed that  $R_i$  is the diagonal matrix,  $\sigma^2 I_{n_i}$ , where  $I_{n_i}$  denotes an  $n_i \times n_i$  identity matrix. In that case,  $\epsilon$

$\epsilon_{ij}$  and  $\epsilon_{ik}$  are uncorrelated, with equal variance, and the  $\epsilon_{ij}$ 's can be thought of as sampling or measurement errors. In principle, we can allow correlation among the  $\epsilon_{ij}$ 's by assuming  $R_i$  has a covariance pattern of the kind considered in Section 7.4. However, doing so would raise two potential complications. First, the  $\epsilon_{ij}$ 's would no longer have a simple interpretation as measurement or sampling errors. This would alter the interpretation of the  $\epsilon_{ij}$ 's, and hence  $b_i$ , implying that the  $\epsilon_{ij}$ 's include a component of model misspecification at the individual level. Second, there can be subtle issues of model identification when  $R_i$  is assumed to have a non-diagonal covariance pattern since there may be insufficient information in the data at hand to support separate estimation of both  $G$  and a non-diagonal  $R_i$ . For example, it is not possible to estimate both  $G$  and an unstructured  $R_i$ . Throughout the remainder of this chapter we assume that the  $\epsilon_{ij}$ 's are pure measurement or sampling errors and that  $R_i = \sigma^2 I_{n_i}$ .

Although we have assumed multivariate normality for both the random effects,  $b_i$ , and the measurement errors,  $\epsilon_i$ , these distributional assumptions are not required for the model development. The form of the conditional and marginal means only requires that the measurement errors are independent of the random effects and that both have mean zero,  $E(b_i) = 0$  and  $E(\epsilon_i) = 0$ . The multivariate normal assumption is required in subsequent sections where we consider estimation, testing, and prediction of random effects.

To clarify the vector and matrix notation introduced so far, consider the following linear mixed effects model with intercepts and slopes that vary randomly among individuals (see [Figure 8.3](#)). For the  $i^{th}$  subject at the  $j^{th}$  measurement occasion, assume that

$$Y_{ij} = \beta_1 + \beta_2 t_{ij} + b_{1i} + b_{2i} t_{ij} + \epsilon_{ij}, \quad j = 1, \dots, n_i.$$

Using vector and matrix notation, this model can be expressed as

$$Y_i = X_i \beta + Z_i b_i + \epsilon_i,$$

where

$$X_i = Z_i = \begin{pmatrix} 1 & t_{i1} \\ 1 & t_{i2} \\ \vdots & \vdots \\ 1 & t_{in_i} \end{pmatrix}.$$

Here  $q = p = 2$  and  $Z_i$  is composed of the two columns of  $X_i$ . This model posits that individuals vary not only in their baseline level of response (when  $t_{i1} = 0$ ), but also in terms of their changes in the mean response over time. The effects of covariates (e.g., due to treatments, exposures, or background characteristics of the individuals) can be incorporated by allowing the means of the intercepts and slopes to depend on these covariates (e.g., by allowing them to vary across the different treatment groups or levels of exposure).

For example, consider the hypothetical two-group study comparing a *treatment* and a *control* group discussed in Section 5.2. If the mean response changes in an approximately linear fashion over time, but with the means of the intercepts and slopes depending on group, the following linear mixed effects model can be adopted:

$$Y_{ij} = \beta_1 + \beta_2 t_{ij} + \beta_3 \text{Group}_i + \beta_4 t_{ij} \times \text{Group}_i + b_{1i} + b_{2i} t_{ij} + \epsilon_{ij},$$

where  $\text{Group}_i = 1$  if the  $i^{th}$  individual was assigned to the treatment, and  $\text{Group}_i = 0$  otherwise. In this model the design matrix  $X_i$  has the following form for the control group:

$$X_i = \begin{pmatrix} 1 & t_{i1} & 0 & 0 \\ 1 & t_{i2} & 0 & 0 \\ \vdots & \vdots & \vdots & \vdots \\ 1 & t_{in_i} & 0 & 0 \end{pmatrix},$$

whereas for the treatment group the design matrix is given by

$$X_i = \begin{pmatrix} 1 & t_{i1} & 1 & t_{i1} \\ 1 & t_{i2} & 1 & t_{i2} \\ \vdots & \vdots & \vdots & \vdots \\ 1 & t_{in_i} & 1 & t_{in_i} \end{pmatrix}.$$

Note that the design matrix  $Z_i$  has the same form for both the treatment and control groups,

$$Z_i = \begin{pmatrix} 1 & t_{i1} \\ 1 & t_{i2} \\ \vdots & \vdots \\ 1 & t_{in_i} \end{pmatrix}.$$

Next consider the covariance among the components of  $Y_i$  in this linear mixed effects model with randomly varying intercepts and slopes. Let  $\text{Var}(b_{1j}) = g_{11}$ ,  $\text{Var}(b_{2i}) = g_{22}$ , and  $\text{Cov}(b_{1i}, b_{2i}) = g_{12}$ . These are the three unique elements of the  $(2 \times 2)$  covariance matrix  $G = \text{Cov}(b_i)$ . If we also assume that  $R_i = \text{Cov}(\epsilon_i) = \sigma^2 I_{n_i}$ , then it can be shown that

$$\begin{aligned} \text{Var}(Y_{ij}) &= \text{Var}(X'_{ij}\beta + Z'_{ij}b_i + \epsilon_{ij}) \\ &= \text{Var}(Z'_{ij}b_i + \epsilon_{ij}) \\ &= \text{Var}(b_{1i} + b_{2i}t_{ij} + \epsilon_{ij}) \\ &= \text{Var}(b_{1i}) + 2t_{ij}\text{Cov}(b_{1i}, b_{2i}) + t_{ij}^2\text{Var}(b_{2i}) + \text{Var}(\epsilon_{ij}) \\ &= g_{11} + 2t_{ij}g_{12} + t_{ij}^2g_{22} + \sigma^2. \end{aligned}$$

Similarly it can be shown that

$$\begin{aligned} \text{Cov}(Y_{ij}, Y_{ik}) &= \text{Cov}(X'_{ij}\beta + Z'_{ij}b_i + \epsilon_{ij}, X'_{ik}\beta + Z'_{ik}b_i + \epsilon_{ik}) \\ &= \text{Cov}(Z'_{ij}b_i + \epsilon_{ij}, Z'_{ik}b_i + \epsilon_{ik}) \\ &= \text{Cov}(b_{1i} + b_{2i}t_{ij} + \epsilon_{ij}, b_{1i} + b_{2i}t_{ik} + \epsilon_{ik}) \\ &= \text{Var}(b_{1i}) + (t_{ij} + t_{ik})\text{Cov}(b_{1i}, b_{2i}) + t_{ij}t_{ik}\text{Var}(b_{2i}) \\ &= g_{11} + (t_{ij} + t_{ik})g_{12} + t_{ij}t_{ik}g_{22}. \end{aligned}$$

Thus in this model for longitudinal data the covariance matrix,  $\text{Cov}(Y_i)$ , can be expressed as a function of time,  $t_{ij}$ . In particular, with the inclusion of random intercepts and slopes, the variance can increase or decrease over time as a quadratic function of the times of measurement. For example, the quadratic expression for  $\text{Var}(Y_{ij})$  given above implies that the variance increases over time (for  $t_{ij} \geq 0$ ) when  $\text{Cov}(b_{1i}, b_{2i}) \geq 0$  but can decrease over time when  $\text{Cov}(b_{1i}, b_{2i}) < 0$ . Similarly the magnitude of the covariance (and correlation) between a pair of responses, say  $Y_{ij}$  and  $Y_{ik}$ , depends on the time separation between them ( $t_{ij}$  and  $t_{ik}$ ). In Section 8.3 we consider the form of the induced random effects covariance structure in the more general case.

Note that the covariance matrix,  $G$ , for the vector of random effects is not invariant to a linear transformation of  $Z_i$ . Linear transformations of the columns of  $Z_i$  alter the interpretation of  $b_i$  and change the estimates of the variances and covariances of the random effects. For example, in the linear mixed effects model with randomly varying intercepts and slopes, centering of the times of measurement (e.g.,  $t_{ij} - \tau$ , for  $\tau \neq 0$ ) alters the interpretation of the intercepts, and this leads to a change in the estimated variance of the random intercepts and their covariance with the random slopes. For example, in the model with untransformed times of measurement the variance of the “intercepts” is a measure of the between-subject variability in the response at time zero. However, in the model with transformed times of measurement, say centered at  $\tau \neq 0$ , the variance of the “intercepts” is a measure of the between-subject variability in the response at time  $\tau$ . Centering not only changes the variance of the random intercepts but also changes the correlation between the random intercepts and slopes. Linear transformations of components of  $Z_i$  produce equivalent mixed effects models only when the covariance matrix,  $G$ , has been left unstructured. When  $G$  is unstructured the appropriate changes to the variances and covariances of the random effects can be produced. For this reason we strongly recommend that the covariance matrix,  $G$ , should always be left unstructured (unless there are compelling reasons, related to the specific analysis under consideration, that suggest this recommendation be relaxed).

Finally, an important issue in the linear mixed effects model concerns the “centering” of the times of measurement. In Chapter 6 we emphasized that “centering” can avoid problems of collinearity when the model for the mean includes linear, quadratic (and possibly higher-order polynomial) time trends. In the linear mixed effects model “centering” has implications for the proper interpretation of

both the mean response and the variance of the random effects. In the illustration above, if  $t_{ij}$  represents time since baseline, then  $\beta_1 + b_{1i}$  represents the subject-specific mean response at baseline (in the control group) and  $\text{Var}(b_{1i}) = g_{11}$  is the between-subject variation in the mean response at baseline. On the other hand, if  $t_{ij}$  is an individual's age at the  $j^{th}$  measurement occasion, then  $\beta_1 + b_{1i}$  does not have a useful interpretation since it represents the subject-specific mean response at age zero; similarly,  $\text{Var}(b_{1i})$  does not have a useful interpretation. In that case there are two obvious choices for centering: (1) center the times of measurement for all subjects at some common fixed age within the age range of the study participants (i.e.,  $t_{ij} - a$ , for some fixed value  $a$ ), or (2) center at the mean age of each subject, when averaged over the subject's period of follow-up (i.e.,  $t_{ij} - \bar{a}_i$ , where  $\bar{a}_i$  is the average age, over the period of follow-up, for the  $i^{th}$  subject). The first option is preferable because  $\beta_1 + b_{1i}$  represents the subject-specific mean response at the common age  $a$  and  $g_{11}$  is the between-subject variation in the mean response at that age. The second option should be avoided because  $\beta_1 + b_{1i}$  then represents the subject-specific mean response at a specific subject's mean age over the period of follow-up. Since the mean age may vary considerably from one subject to another,  $g_{11}$  will be inflated and will not have a meaningful interpretation. In summary, with unbalanced longitudinal data, mean centering of the times of measurement should be avoided. Instead, we recommend that times of measurement should be centered at some common value of time (or age) in the center of the range of values for all individuals. By centering at a common value, the intercept is interpretable as the mean response at that common value for time (or age) and  $\text{Var}(b_{1i})$  also has a meaningful interpretation.

## 8.3 RANDOM EFFECTS COVARIANCE STRUCTURE

Next we consider the form of the induced random effects covariance structure for longitudinal data in the more general case. In the linear mixed effects model

$$Y_i = X_i\beta + Z_i b_i + \epsilon_i,$$

$R_i = \text{Cov}(\epsilon_i)$  describes the covariance among the longitudinal observations when focusing on the conditional mean response profile of a *specific* individual. That is, it is the covariance of the  $i^{\text{th}}$  individual's deviations from her mean response profile,

$$E(Y_i|b_i) = X_i\beta + Z_i b_i.$$

For example, in [Figures 8.2](#) and [8.3](#) these deviations are positive and negative, and vary randomly about zero, for individuals A and B. As mentioned previously, it is usually assumed that  $R_i$  is a diagonal matrix,  $\sigma^2 I_{n_i}$ , where  $I_{n_i}$  denotes an  $n_i \times n_i$  identity matrix. The latter is often referred to as a “conditional independence assumption”; that is, given the random effects  $b_i$ , the measurement errors are independently distributed with a common variance  $\sigma^2$ .

In the linear mixed effects model we can distinguish the conditional mean of  $Y_i$ , given  $b_i$ ,

$$E(Y_i|b_i) = X_i\beta + Z_i b_i,$$

from the *marginal* or population-averaged mean of  $Y_i$ ,

$$E(Y_i) = X_i\beta,$$

where averaging is over the distribution of the random effects,  $b_i$ . In a similar way we can distinguish between conditional and marginal covariances. The conditional covariance of  $Y_i$ , given  $b_i$ , is

$$\text{Cov}(Y_i|b_i) = \text{Cov}(\epsilon_i) = R_i,$$

while the marginal covariance of  $Y_i$ , averaged over the distribution of  $b_i$ , is

$$\begin{aligned} \text{Cov}(Y_i) &= \text{Cov}(Z_i b_i) + \text{Cov}(\epsilon_i) \\ &= Z_i \text{Cov}(b_i) Z'_i + \text{Cov}(\epsilon_i) \\ &= Z_i G Z'_i + R_i. \end{aligned}$$

This latter expression for the marginal covariance may be somewhat daunting at first glance. Even when  $R_i = \text{Cov}(\epsilon_i) = \sigma^2 I_{n_i}$ , a diagonal matrix (with all pairwise correlations equal to zero),

$$\text{Cov}(Y_i) = Z_i G Z'_i + \sigma^2 I_{n_i}$$

is emphatically not a diagonal matrix. That is,  $\text{Cov}(Y_i)$  will, in general, have non-zero off-diagonal elements, thereby accounting for the correlation among the repeated observations on the same individuals in a longitudinal study. Thus the introduction of random effects,  $b_i$ , induces correlation among the components of  $Y_i$ . An additional property of the linear mixed effects model is that  $\text{Cov}(Y_i)$  has been described in terms of a set of covariance parameters, some defining the matrix  $G$  and some defining the matrix  $R_i$ . That is, the linear mixed effects model allows for the explicit analysis of between-subject ( $G$ ) and within-subject ( $R_i$ ) sources of variation in the responses. Finally, the marginal covariance of  $Y_i$  is a function of the times of measurement. For example, in the model with randomly varying intercepts and slopes considered in Section 8.2, we saw that

$$\text{Var}(Y_{ij}) = g_{11} + 2t_{ij}g_{12} + t_{ij}^2g_{22} + \sigma^2,$$

$$\text{Cov}(Y_{ij}, Y_{ik}) = g_{11} + (t_{ij} + t_{ik})g_{12} + t_{ij}t_{ik}g_{22},$$

so that both  $\text{Var}(Y_{ij})$  and  $\text{Cov}(Y_{ij}, Y_{ik})$  depend on the measurement times.

The induced random effects covariance structure,

$$\text{Cov}(Y_i) = Z_i G Z'_i + \sigma^2 I_{n_i},$$

can be contrasted with the covariance pattern models described in Chapter 7 (see Section 7.4). Recall that a defining feature of the covariance pattern models is that they take into account all sources of variability that have an impact on the covariance, but do not distinguish between the different sources of variability. In contrast, linear mixed effects models explicitly distinguish

between-subject and within-subject sources of variability. For linear models for longitudinal continuous data, both approaches yield the same model for the *marginal* or population-averaged mean of  $Y_i$ ,

$$E(Y_i) = X_i\beta,$$

and differ only in terms of the assumed model for the covariance. As we will see in Chapters 12 through 16, longitudinal models for discrete responses do not share this property; for discrete responses, different approaches for accounting for the covariance among the longitudinal responses can lead to models for the mean response having regression parameters with quite distinct interpretations.

The induced random effects covariance structure has certain features that are different from the covariance pattern models considered in Chapter 7. First, unlike many covariance pattern models, the random effects covariance structure does not require a balanced longitudinal design. Because the covariance is expressed as an explicit function of the times of measurement (when times of measurement, or functions of time, are included in  $Z_i$ ), in principle, each individual can have a unique sequence of measurement times. This makes linear mixed effects models well suited for modeling data from inherently unbalanced longitudinal designs. In addition the number of covariance parameters is the same regardless of the number and timing of the measurements. Finally, unlike many of the covariance pattern models that make strong assumptions about homogeneity of variance over time, the random effects covariance structure allows the variance and covariance to increase or decrease as a function of the times of measurement (e.g., in the random intercepts and slopes model, the variance is a quadratic function of the times of measurement).

## 8.4 TWO-STAGE RANDOM EFFECTS FORMULATION

The linear mixed effects model given by (8.2) can be motivated by a two-stage random effects formulation of the model. Indeed, some of the main ideas behind the mixed effects model are often better understood by considering the model as arising from a two-stage specification. For purely pedagogical purposes, we find the two-stage specification to be quite helpful; however, we must caution the reader that the two-stage formulation of the linear mixed effects model does introduce some unnecessary restrictions on the model.

# Stage 1

As the term implies, a two-stage random effects model can be conceived in two separate stages. In the first stage subjects are assumed to have their own unique individual-specific mean response trajectories. That is, in stage 1 we posit that the repeated measures on each individual follow a regression model having the same set of covariates, but with separate or distinct regression coefficients for each individual. This is expressed more formally as

$$Y_i = Z_i \beta_i + \epsilon_i,$$

where the vector of errors,  $\epsilon_i$ , are assumed to have a normal distribution, with mean equal to zero and variance  $\sigma^2$ . That is, the  $\epsilon_i$  can be thought of as measurement or sampling errors, with  $\epsilon_i \sim N(0, \sigma^2 I_{n_i})$ . Note that the number of individual-specific regression coefficients is the same (i.e., the dimension of  $\beta_i$  is  $q$ ), regardless of the number of longitudinal responses  $n_i$ . These individual-specific regression coefficients,  $\beta_i$ , can be interpreted as the  $i^{th}$  individual's "true" regression coefficients. Alternatively,  $Z_i \beta_i$  can be thought of as the  $i^{th}$  individual's "true" underlying mean response trajectory. When viewed in this way, the longitudinal responses on the  $i^{th}$  individual are assumed to follow the individual-specific response trajectory given by  $Z_i \beta_i$ , but with the addition of measurement or sampling errors,  $\epsilon_i$ .

Note that the matrix  $Z_i$  specifies how an individual's mean response changes over time and/or how the mean response changes with other time-varying covariates (e.g., height). For example, it might be assumed that the mean response trajectory is linear, quadratic, or a spline function of time. Consider a model that assumes the individual-specific trajectories are linear in time. Then, the first-stage model can be written as

$$\begin{pmatrix} Y_{i1} \\ Y_{i2} \\ \vdots \\ Y_{in_i} \end{pmatrix} = \begin{pmatrix} 1 & t_{i1} \\ 1 & t_{i2} \\ \vdots & \vdots \\ 1 & t_{in_i} \end{pmatrix} \begin{pmatrix} \beta_{1i} \\ \beta_{2i} \end{pmatrix} + \begin{pmatrix} \epsilon_{i1} \\ \epsilon_{i2} \\ \vdots \\ \epsilon_{in_i} \end{pmatrix}.$$

The essential idea underlying the first-stage model is to fit separate linear regression models to the data for each individual, but with the proviso that these regressions involve the same set of covariates,  $Z_i$ . This is an important observation since it implies that, in principle (and given a sufficient number of repeated measures on each individual), it should be possible to estimate  $\beta_i$  (and  $\sigma^2$ ) using data from only the  $i^{th}$  individual.

Recall that a particular feature of this first-stage formulation is that the matrix of covariates  $Z_i$  is restricted to contain only within-individual or time-varying covariates (with the exception of the column of 1's for the intercept). Time-invariant or between-individual covariates (e.g., gender, treatment group, exposure group) cannot be included in  $Z_i$  since their effects would simply be absorbed into the intercept term. Instead, between-individual covariates are introduced in the second stage of the model formulation.

## Stage 2

In the second stage we make the assumption that the individual-specific effects,  $\beta_i$ , are random. Given that the  $\beta_i$  are random variables, they have some probability distribution, with a mean and covariance. The mean and covariance of the  $\beta_i$  are the population parameters that are modeled in the second stage. Specifically, variation in  $\beta_i$  from one individual to another is modeled as a function of a set of between-individual (or time-invariant) covariates (e.g., gender, treatment group). In particular, the mean of the  $\beta_i$  can be expressed as a linear function of a set of between-individual covariates,  $A_i$ ,

$$E(\beta_i) = A_i\beta,$$

where  $A_i$  is a  $q \times p$  matrix. The remaining residual between-individual variation in the  $\beta_i$  that cannot be explained by  $A_i$  is expressed as

$$\text{Cov}(\beta_i) = G.$$

Specification of a model for the mean and covariance of the  $\beta_i$  completes the second stage of the model formulation.<sup>1</sup>

For example, consider the hypothetical two-group study comparing a *treatment* and a *control* discussed earlier. If we assume that individual-specific changes in the mean response over time are linear, the first-stage model is given by

$$\begin{pmatrix} Y_{i1} \\ Y_{i2} \\ \vdots \\ Y_{in_i} \end{pmatrix} = \begin{pmatrix} 1 & t_{i1} \\ 1 & t_{i2} \\ \vdots & \vdots \\ 1 & t_{in_i} \end{pmatrix} \begin{pmatrix} \beta_{1i} \\ \beta_{2i} \end{pmatrix} + \begin{pmatrix} \epsilon_{i1} \\ \epsilon_{i2} \\ \vdots \\ \epsilon_{in_i} \end{pmatrix}.$$

In the second stage, we can allow the mean of  $\beta_i$  (i.e., the mean intercept and slope) to depend on group. For example, a model that allows both the mean intercept and slope to depend on group is given by

$$E(\beta_{1i}) = \beta_1 + \beta_3 \text{Group}_i$$

$$E(\beta_{2i}) = \beta_2 + \beta_4 \text{Group}_i$$

where  $\text{Group}_i = 1$  if the  $i^{th}$  individual was assigned to the treatment, and  $\text{Group}_i = 0$  otherwise. In this model,  $\beta_1$  is the mean intercept in the control group, while  $\beta_1 + \beta_3$  is the mean intercept in the treatment group. That is,  $\beta_3$  represents the treatment group difference in the mean intercept. When  $t_{ij}$  is the time since baseline,  $\beta_3$  has a useful interpretation in terms of a treatment group difference in the mean response at baseline. Similarly  $\beta_2$  is the mean slope, or rate of change in the mean response over time, in the control group, while  $\beta_2 + \beta_4$  is the mean slope in the treatment group. That is,  $\beta_4$  has interpretation in terms of a treatment group difference in the mean slope or rate of change in the mean response over time. In this model the design matrix  $A_i$  of between-individual covariates has the following form:

$$A_i = \begin{pmatrix} 1 & 0 & \text{Group}_i & 0 \\ 0 & 1 & 0 & \text{Group}_i \end{pmatrix}.$$

Thus, for the control group, the model for the mean is

$$E\left(\begin{pmatrix} \beta_{1i} \\ \beta_{2i} \end{pmatrix}\right) = \begin{pmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \end{pmatrix} \begin{pmatrix} \beta_1 \\ \beta_2 \\ \beta_3 \\ \beta_4 \end{pmatrix} = \begin{pmatrix} \beta_1 \\ \beta_2 \end{pmatrix};$$

similarly, for the treatment group, the model for the mean is

$$E\left(\begin{pmatrix} \beta_{1i} \\ \beta_{2i} \end{pmatrix}\right) = \begin{pmatrix} 1 & 0 & 1 & 0 \\ 0 & 1 & 0 & 1 \end{pmatrix} \begin{pmatrix} \beta_1 \\ \beta_2 \\ \beta_3 \\ \beta_4 \end{pmatrix} = \begin{pmatrix} \beta_1 + \beta_3 \\ \beta_2 + \beta_4 \end{pmatrix}.$$

It is also assumed that the remaining residual variation in  $\beta_i$ , which cannot be explained by the effect of group, is

$$\text{Cov}(\beta_i) = G = \begin{pmatrix} g_{11} & g_{12} \\ g_{21} & g_{22} \end{pmatrix},$$

where  $g_{11} = \text{Var}(\beta_{1i})$ ,  $g_{22} = \text{Var}(\beta_{2i})$ , and  $g_{12} = g_{21} = \text{Cov}(\beta_{1i}, \beta_{2i})$ . Thus  $g_{11}$  is the variance of  $\beta_{1i}$ , after adjusting for the effect of treatment group, and so on.

The two components of the two-stage model can be combined to yield a linear mixed effects model for  $Y_i$ , albeit one that has some restrictions. To see how this can be achieved, let us rewrite the subject-specific effects,  $\beta_i$ , as

$$\beta_i = A_i\beta + b_i,$$

where  $b_i$  has a multivariate normal distribution with mean zero and covariance matrix,  $G$ . Here the  $b_i$  yield the regression coefficients from an individual's *residual* trajectory over time, after the covariate effects have been accounted for. Put another way, the  $b_i$  represent the  $i^{\text{th}}$  individual's deviation from the population mean response. Next, by combining the two components of the two-stage model, we obtain

$$\begin{aligned} Y_i &= Z_i\beta_i + \epsilon_i \\ &= Z_i(A_i\beta + b_i) + \epsilon_i \\ &= (Z_iA_i)\beta + Z_ib_i + \epsilon_i \\ &= X_i\beta + Z_ib_i + \epsilon_i, \end{aligned}$$

where  $X_i = Z_iA_i$ . When averaged over the random effects,  $b_i$ ,

$$E(Y_i) = (Z_iA_i)\beta = X_i\beta,$$

and

$$\text{Cov}(Y_i) = Z_iGZ'_i + \sigma^2 I_{n_i}.$$

Note that in the two-stage formulation,  $Z_i$  appears in both the models for the marginal mean and covariance.

While this model is remarkably similar to the linear mixed effects model introduced in the previous section, there is one important difference. The two-stage model places a constraint on the choice of the design matrix for the fixed effects. That is, the two-stage formulation requires that the design matrix for the fixed effects has the special structure  $X_i = Z_iA_i$ , where  $A_i$  contains only between-subject (or time-invariant) covariates and  $Z_i$  contains only within-subject (or time-varying) covariates. This form for the design matrix for the fixed effects implies that any time-varying covariates must be specified as random effects to ensure their inclusion in the model for the population mean response.<sup>2</sup> This constraint is unnecessary and, in many settings, it can be somewhat inconvenient. In some applications this constraint forces us to consider rather more complex models than may be necessary. For example, in order to allow a sufficiently complex model for the mean response over time (specified in terms of  $Z_iA_i\beta$ ), it may be necessary to include many covariates in  $Z_i$ . However, in the two-stage model formulation, this can only be achieved by also introducing an equally complex model for the covariance, since

$$\text{Cov}(Y_i) = Z_iGZ'_i + \sigma^2 I_{n_i}.$$

An example of this arises in developing a model for FEV<sub>1</sub> in children. Previous studies have shown that both age (as a linear spline) and log height are important covariates. Thus four subject-specific regression coefficients (an intercept, two coefficients for age, and one coefficient for log height) are needed to model the mean. But a  $4 \times 4$  covariance matrix for  $G$  is very unwieldy and difficult to fit without very large samples.

Alternatively, in the two-stage formulation a very simple structure for the covariance imposes an often unrealistically simple structure on the mean response. The most extreme example of this is the two-stage model, which induces a compound symmetry covariance. In that case a compound symmetry covariance is obtained from a two-stage model with randomly varying intercepts,

$$Y_i = Z_i\beta_i + \epsilon_i,$$

where  $Z_i$  is a  $(n_i \times 1)$  vector of 1's. While marginally (or averaged over the random effects) this model induces a simple compound symmetry covariance structure

$$\text{Cov}(Y_i) = Z_iGZ'_i + \sigma^2 I_{n_i} = g_{11}J_{n_i} + \sigma^2 I_{n_i},$$

where  $J_{n_i}$  denotes a  $(n_i \times n_i)$  matrix of 1's, this model precludes any dependence of the mean response on time. That is,

$$E(Y_i) = (Z_i A_i)\beta$$

cannot depend on time since time, a within-subject covariate, has not been included in  $Z_i$  in the first stage. Thus, when formulated as a two-stage model, the randomly varying intercepts model excludes the most salient within-subject covariate in a longitudinal study (i.e., time), and thereby does not allow for estimation of changes in the mean response over time, the primary goal in a longitudinal analysis!

In summary, we view the two-stage formulation as being most useful for motivating the main ideas and concepts underlying linear mixed effects models. The inherent restrictions in the two-stage formulations can be circumvented by considering linear mixed effects models with an arbitrary design matrix,  $X_i$ , for the fixed effects, and by allowing the dimension of  $Z_i$  to be arbitrary. For many models for longitudinal analysis the only restriction placed on  $X_i$  and  $Z_i$  is that  $Z_i$  is composed of a subset of the columns of  $X_i$  (in Chapter 19 we will encounter linear mixed effects models where the columns of  $Z_i$  are no longer a subset of the columns of  $X_i$ ). The latter constraint ensures that  $Z_i b_i$  can be interpreted as a zero mean between-subject residual trajectory (or, put another way, the discrepancy between the  $i^{th}$  individual's conditional mean response trajectory and the mean response trajectory in the population). Thus in the linear mixed effects model,

$$Y_i = X_i \beta + Z_i b_i + \epsilon_i,$$

the restriction that the columns of  $Z_i$  are a subset of the columns of  $X_i$  allows us to partition the columns of  $X_i$  into a set of columns corresponding to the effects that are fixed and a complementary set of columns corresponding to the effects that are random. If we denote the former by  $X_i^{(F)}$  and the latter by  $X_i^{(R)}$ , the model for  $Y_i$  can then be rewritten as

$$Y_i = X_i^{(F)} \beta^{(F)} + X_i^{(R)} \beta_i^{(R)} + \epsilon_i,$$

awhere  $\beta$  has been similarly partitioned into effects that are considered to be fixed,  $\beta^{(F)}$ , and effects that are considered to be random,  $\beta_i^{(R)}$ .

Finally, on an historical note, a version of the two-stage formulation was popularized by biostatisticians working at the U.S. National Institutes of Health (NIH). They proposed a method for analyzing repeated measures data where in the first stage subject-specific regression coefficients are estimated using ordinary least-squares regression (based only on the observations for each subject). In the second stage, the estimated regression coefficients are then analyzed as summary measures using standard parametric (or nonparametric) methods. This method for analyzing repeated measures data became known as the “NIH method”<sup>3</sup> and is a variant of the summary measure analyses considered in Section 3.6.

## 8.5 CHOICE AMONG RANDOM EFFECTS COVARIANCE MODELS

Although the linear mixed effects model assumes that the longitudinal responses depend on a combination of population and subject-specific effects, when averaged over the distribution of the random effects,

$$E(Y_i) = X_i\beta,$$

and the covariance among the responses has the distinctive random effects structure,

$$\text{Cov}(Y_i) = Z_i G Z_i' + \sigma^2 I_{n_i}.$$

From the perspective of modeling the covariance, the random effects structure is appealing because the number of covariance parameters,  $q \times (q + 1)/2 + 1$ , is the same regardless of the number and timing of the measurement occasions. In many applications, it will be sufficient to include only random intercepts and slopes for time (a total of  $2 \times (2 + 1)/2 + 1 = 4$  covariance parameters), thereby allowing for heterogeneity in the variances and correlations that can be expressed as functions of time. In other applications, a more complex random effects structure may be required.

In choosing a model for the covariance, it will often be of interest to compare two nested models, one with  $q$  correlated random effects, the other with  $q + 1$  correlated random effects. The difference in the number of covariance parameters between these two models is  $q + 1$ , since the “full” model contains one additional variance and  $q$  additional covariances. As mentioned in Section 7.5, in general, the likelihood ratio test provides a valid method for comparing nested models for the covariance. However, in certain cases, the usual null distribution for the likelihood ratio test is no longer valid. In particular, the comparison of random effects models for the covariance is such a non-standard problem.

In general, when testing a null hypothesis that is on the boundary of the parameter space (e.g., the variance of a random effect equals zero), the usual null distribution for the likelihood ratio test is no longer a chi-squared distribution with degrees of freedom equal to the difference between the number of parameters in the full and reduced models; instead, the null distribution is a mixture of chi-squared distributions. For example, when comparing two nested models, one with  $q$  correlated random effects, the other with  $q + 1$  correlated random effects, the null distribution of the likelihood ratio test is a 50:50 mixture of chi-squared distributions with  $q$  and  $q + 1$  degrees of freedom. In Section 7.5 we considered the special case where  $q = 0$ . A table of critical values, when the null distribution is a known 50:50 mixture of chi-squared distributions with  $q$  and  $q + 1$  degrees of freedom, is provided in [Table C.1](#) in Appendix C. The critical values in [Table C.1](#) can be used for making inferences about the complexity of the random effects covariance structure. For example, when comparing a model with two correlated random effects (e.g., random intercepts and slopes) versus a model with one random effect (e.g., random intercepts only), the critical value for the 0.05 significance level can be found in the second row ( $q = 1$ ) of [Table C.1](#). This yields a critical value of 5.14, which is somewhat smaller than the critical value of 5.99 from a standard chi-squared distribution with 2 degrees of freedom.

Alternatively, and especially for more complex comparisons among nested random effects models for the covariance where the null distribution of the likelihood ratio test is not well understood (e.g., comparisons of nested models with  $q$  correlated random effects and  $q + k$  correlated random effects, where  $k > 1$ ), we recommend the use of  $\alpha = 0.1$ , instead of  $\alpha = 0.05$ , when judging the statistical significance of the likelihood ratio test. The latter procedure is somewhat ad hoc but will protect against selection of a model that is too parsimonious. In conclusion, for simple comparisons among nested random effects models, the likelihood ratio test statistic can be compared with the critical values in [Table C.1](#). For more complex comparisons, we recommend the use of the  $\alpha = 0.1$ , instead of  $\alpha = 0.05$ , significance level.

## 8.6 PREDICTION OF RANDOM EFFECTS

In this section we provide a non-technical discussion on the prediction of random effects. A good grasp of the material in this section is all that is required for an understanding of the notion of predicting random effects. In Section 8.7 we present a more detailed and technical discussion of the same topic. Many of our readers may find the level of mathematical difficulty of the material in Section 8.7 too challenging. While we encourage all of our readers to tackle Section 8.7, we note that it can be omitted at first reading without loss of continuity.

In many applications where longitudinal data arise, inference is focused on the fixed effects,  $\beta$ . These regression parameters have interpretation in terms of changes in the mean response over time, and their relation to covariates. However, in some longitudinal studies, we may want to predict subject-specific response profiles. For example, in studies of growth it may be of interest to obtain subject-specific growth trajectories. In other types of longitudinal studies, it may be of interest to identify those individuals who showed the greatest increase or decrease in the response over time. Because the linear mixed effects model explicitly distinguishes between fixed and random effects, we can also estimate (or predict) individual-specific response trajectories over time. That is, it is possible to obtain predictions of the subject-specific effects,  $b_i$ , or of the subject-specific response trajectories,  $X_i\beta + Z_i b_i$ . Technically, because the  $b_i$  are random variables, and not fixed population parameters, we customarily refer to “predicting” the random effects rather than “estimating” them.

In general, the problem of predicting a random variable can be shown to be that of predicting its conditional mean, given the available data. Thus the best predictor of  $b_i$  is the conditional mean of  $b_i$ , given the vector of responses  $Y_i$  (and  $\hat{\beta}$ ),

$$E(b_i|Y_i) = GZ_i'\Sigma_i^{-1}(Y_i - X_i\hat{\beta}),$$

where  $\sum_i = \text{Cov}(Y_i) = Z_i G Z_i' + R_i$ . This is known as the “best linear unbiased predictor” (or BLUP). This predictor of  $b_i$  depends on the unknown covariance among the  $Y_i$ . When the unknown covariance parameters are replaced by their REML (or ML) estimates, the resulting predictor

$$\hat{b}_i = \hat{G}Z_i'\hat{\Sigma}_i^{-1}(Y_i - X_i\hat{\beta}),$$

is referred to as the “empirical BLUP.” Given the “empirical BLUP,”  $\hat{b}_i$ , we can also obtain the  $i^{th}$  subject’s predicted response profile as follows:

$$\hat{Y}_i = X_i\hat{\beta} + Z_i\hat{b}_i.$$

Interestingly, the  $i^{th}$  subject’s predicted response profile can also be expressed as a weighted average of the estimated population-averaged mean response profile,  $X_i\hat{\beta}$ , and the  $i^{th}$  subject’s observed response profile  $Y_i$ . Specifically,

$$\hat{Y}_i = X_i\hat{\beta} + Z_i\hat{b}_i = (\hat{R}_i\hat{\Sigma}_i^{-1})X_i\hat{\beta} + (I_{n_i} - \hat{R}_i\hat{\Sigma}_i^{-1})Y_i.$$

This expression helps to explain why it is often said that the empirical BLUP estimator “shrinks” the  $i^{th}$  subject’s predicted response profile towards the population-averaged mean response profile.

The amount of “shrinkage” toward the population mean depends on the relative magnitude of  $R_i$  and  $\sum_i$ . Recall that  $R_i$  characterizes the within-subject variability, while  $\sum_i$  incorporates both within-subject and between-subject sources of variability. As a result, when  $R_i$  is relatively “large,” and the within-subject variability is greater than the between-subject variability, more weight is assigned to  $X_i\hat{\beta}$ , the estimated population-averaged mean response profile, than to the  $i^{th}$  individual’s observed responses. On the other hand, when the between-subject variability is large relative to the within-subject variability, more weight is given to the  $i^{th}$  subject’s observed responses,  $Y_i$ . Intuitively, this weighting scheme is quite sensible since greater weight should be given to the  $i^{th}$  individual’s observed responses when any within-subject variability in the longitudinal responses (e.g., due to measurement error) is relatively small when compared to the natural heterogeneity in the individual-specific longitudinal response trajectories. On the other hand, less weight should be given to the  $i^{th}$  individual’s observed responses when the within-subject variability is relatively large and the

population is relatively homogeneous. Finally, the amount of “shrinkage” toward the population mean depends also on  $n_i$ , the number of observation on the  $i^{th}$  subject. In general, there is more shrinkage toward the population mean curve when  $n_i$  is small. Intuitively, this is also quite sensible since less weight should be given to the  $i^{th}$  individual’s observed responses when fewer data points are available.

## 8.7 PREDICTION AND SHRINKAGE\*

In this section<sup>†</sup> we present a more detailed and technical discussion on prediction of random effects in the linear mixed effects model. In doing so, we provide some motivation for, and expressions that support, the main results outlined in Section 8.6.

Because the linear mixed effects model explicitly distinguishes between fixed and random effects, we can estimate both types of effects. As noted in the previous section, the prediction of random effects translates into the problem of predicting the conditional mean of  $b_i$ , given the vector of responses  $Y_i$ ,  $E(b_i|Y_i)$ . Under the assumptions of the linear mixed effects model,  $Y_i$  and  $b_i$  have a joint multivariate normal distribution. Using well-known properties of the multivariate normal distribution, it can be shown that the conditional mean of  $b_i$  given  $Y_i$  (and  $\hat{\beta}$ ) is

$$E(b_i|Y_i) = GZ'_i\Sigma_i^{-1}(Y_i - X_i\hat{\beta}),$$

where  $\Sigma_i = \text{Cov}(Y_i) = Z_iGZ'_i + R_i$ . This is known as the “best linear unbiased predictor” (or BLUP). From a practical standpoint, this predictor of  $b_i$  is unusable because it depends on the unknown covariance parameters. When the unknown covariance parameters are replaced by their REML estimates, the resulting predictor

$$\hat{b}_i = \hat{G}Z'_i\hat{\Sigma}_i^{-1}(Y_i - X_i\hat{\beta}),$$

is referred to as the “empirical BLUP” or the “empirical Bayes” (EB) estimator (since  $\hat{b}_i$  can also be derived from a fully Bayesian formulation). In addition to obtaining a prediction of  $b_i$ , we can obtain standard errors for the prediction based on the expression

$$\text{Cov}(\hat{b}_i - b_i) = G - GZ'_i\Sigma_i^{-1}Z_iG + GZ'_i\Sigma_i^{-1}X_i \left( \sum_{i=1}^N X'_i\Sigma_i^{-1}X_i \right)^{-1} X'_i\Sigma_i^{-1}Z_iG.$$

Note that  $\text{Cov}(\hat{b}_i - b_i)$  is used to assess the precision of the prediction of  $b_i$ , rather than  $\text{Cov}(\hat{b}_i)$ , because the latter would fail to recognize the random variation of  $b_i$ .

Standard errors for the prediction are obtained by simply substituting  $\hat{\Sigma}_i$  and  $\hat{G}$ , the REML estimates of the covariance parameters, in the previous expression for  $\text{Cov}(\hat{b}_i - b_i)$ .

Given the prediction of  $b_i$ , the  $i^{th}$  subject’s predicted response profile is given by

$$\hat{Y}_i = X_i\hat{\beta} + Z_i\hat{b}_i.$$

This expression for the predicted response profile can be re-expressed as follows:

$$\begin{aligned} \hat{Y}_i &= X_i\hat{\beta} + Z_i\hat{b}_i \\ &= X_i\hat{\beta} + Z_i\hat{G}Z'_i\hat{\Sigma}_i^{-1}(Y_i - X_i\hat{\beta}) \\ &= (I_{n_i} - Z_i\hat{G}Z'_i\hat{\Sigma}_i^{-1})X_i\hat{\beta} + Z_i\hat{G}Z'_i\hat{\Sigma}_i^{-1}Y_i \\ &= (\hat{R}_i\hat{\Sigma}_i^{-1})X_i\hat{\beta} + (I_{n_i} - \hat{R}_i\hat{\Sigma}_i^{-1})Y_i, \end{aligned}$$

since

$$\hat{\Sigma}_i\hat{\Sigma}_i^{-1} = I_{n_i} = (Z_i\hat{G}Z'_i + \hat{R}_i)\hat{\Sigma}_i^{-1} = Z_i\hat{G}Z'_i\hat{\Sigma}_i^{-1} + \hat{R}_i\hat{\Sigma}_i^{-1}.$$

The latter expression for  $\hat{Y}_i$  shows how the empirical Bayes estimator “shrinks” the  $i^{th}$  subject’s predicted response profile toward the population-averaged mean response profile. As noted in Section 8.6, the amount of “shrinkage” depends on the relative magnitude of  $R_i$  and  $\Sigma_i$ . When the within-subject variability,  $R_i$ , is large relative to the between-subject variability, more weight is assigned to  $X_i\hat{\beta}$  than to the  $i^{th}$  individual’s observed responses. Conversely, when the between-subject variability is large relative to the within-subject variability, more weight is given to the  $i^{th}$  subject’s observed responses,  $Y_i$ .

Similarly it can be shown that the prediction of individual-specific regression coefficients is a weighted average of the REML estimate of the fixed effects and the corresponding OLS estimate based only on the individual’s observations. Specifically, when the linear mixed effects model has a two-stage representation, with  $X_i = Z_iA_i$  and  $R_i = \sigma^2I_{n_i}$ , the “empirical BLUP” of  $\beta_i - A_i\beta + b_i$  is a weighted average of  $A_i\hat{\beta}$  and  $\hat{\beta}_i^{\text{ols}}$ , where  $\hat{\beta}$  is the usual REML estimate of  $\beta$  obtained from the

available data on all subjects and  $\hat{\beta}_i^{\text{OLS}}$  is the ordinary least squares estimate of  $\beta_i$  based only on the  $n_i$  observations for the  $i^{\text{th}}$  subject. More formally, the “empirical BLUP” of  $\beta_i$  can be expressed as

$$\hat{\beta}_i = A_i \hat{\beta} + \hat{b}_i = W_i \hat{\beta}_i^{\text{OLS}} + (I_q - W_i) A_i \hat{\beta},$$

where the “weight,”  $W_i$ , is a ratio of the between-subject variability to the sum of the between- and within-subject variability,

$$W_i = G\{G + \sigma^2(Z_i' Z_i)^{-1}\}^{-1},$$

and  $I_q$  denotes a  $q \times q$  identity matrix. Although this expression for the “weight”,  $W_i$ , appears somewhat daunting, note that when there is very little within-subject variability (and  $\sigma^2 \approx 0$ ),  $W_i \approx I_q$  and then  $\hat{\beta}_i \approx \hat{\beta}_i^{\text{OLS}}$ . That is, when there is very little within-subject variability, we have almost perfect information about  $b_i$  from  $Y_i$  alone. On the other hand, when there is very little between-subject variability (and  $G \approx 0$ ),  $W_i \approx 0$  and then  $\hat{\beta}_i \approx A_i \hat{\beta}$ . Thus, when there is very little between-subject variability in the individual-specific trajectories, it is quite sensible to base our “estimate” or prediction of  $b_i$  on data from all of the individuals in the study. The expression for the weight  $W_i$  also highlights how the number of repeated measurements influences the compromise between  $\hat{\beta}_i^{\text{OLS}}$  and  $A_i \hat{\beta}$ . For example, consider the special case of the model with randomly varying intercepts, where  $Z_i$  is an  $n_i \times 1$  vector of 1’s. It can be shown that

$$W_i = G\{G + \sigma^2(Z_i' Z_i)^{-1}\}^{-1} = \frac{n_i g_{11}}{n_i g_{11} + \sigma^2},$$

where  $g_{11} = \text{Var}(b_{1i})$  is the variance of the random intercept. Thus, for fixed values of the within- and between-subject variability, the more observations that are available on the  $i^{\text{th}}$  individual the more the “empirical BLUP” of  $\beta_i$  relies on that individual’s data to “estimate” the random effect.

## 8.8 CASE STUDIES

Next we illustrate the main ideas presented in this chapter by considering linear mixed effects models for analyzing data from three different studies. The first illustration uses lung function growth data in a sample of children and adolescents from the Six Cities Study of Air Pollution and Health. The second illustration uses data on body fat accretion from a prospective study of the development of obesity in a cohort of girls. The third illustration uses data on CD4 counts from a randomized clinical trial of AIDS patients with advanced immune suppression.

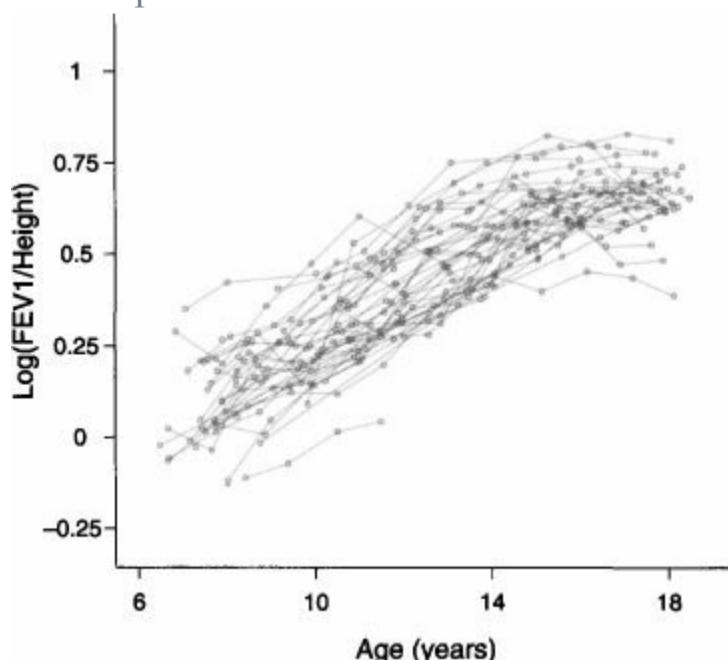
# Six Cities Study of Air Pollution and Health

The Six Cities Study of Air Pollution and Health was a longitudinal study designed to characterize lung growth as measured by changes in pulmonary function in children and adolescents, and the factors that influence lung function growth (Dockery et al., 1983). A cohort of 13,379 children born in or after 1967 was enrolled in six communities across the United States: Watertown (Massachusetts), Kingston and Harriman (Tennessee), a section of St. Louis (Missouri), Steubenville (Ohio), Portage (Wisconsin), and Topeka (Kansas). Most children were enrolled in the first or second grade (between the ages of 6 and 7) and measurements of study participants were obtained annually until graduation from high school or loss to follow-up. At each annual examination, spirometry, the measurement of pulmonary function, was performed and a respiratory health questionnaire was completed by a parent or guardian. The basic maneuver in simple spirometry is maximal inspiration (or breathing in) followed by forced exhalation as rapidly as possible into a closed container. Many different measures can be derived from the spirometric curve of volume exhaled versus time.

One widely used measure is the total volume of air exhaled in the first second of the maneuver ( $FEV_1$ ).

In this section we present an analysis of a subset of the pulmonary function data collected in the Six Cities Study. The data consist of all measurements of  $FEV_1$ , height and age obtained from a randomly selected subset of the female participants living in Topeka, Kansas. The random sample consists of 300 girls, with a minimum of one and a maximum of twelve observations over time; data for four selected girls are presented in [Table 8.1](#). Note that although girls with only a single observation do not directly provide information about longitudinal or intra-individual change over time, their observations at a single occasion do contribute to the analysis (e.g., these observations contribute information to the estimation of variances and between-subject effects). Examination of [Table 8.1](#) reveals that these data are inherently unbalanced over time, and the degree of imbalance is even more marked when the age of the child is used as the metamer for time. That is, in this data set children enter the study at different ages and also have different occasions of measurement. [Figure 8.4](#) displays a time plot, with joined line segments, of  $\log(FEV_1/\text{height})$  versus age for 50 randomly selected girls.

**Fig. 8.4** Time plot, with joined line segments, of  $\log(FEV_1/\text{height})$  versus age in years for 50 randomly selected girls from the Topeka data set.



**Table 8.1** Data on age, height, and  $FEV_1$  for four girls selected from the Topeka data set.

Subject ID	Age	Height	Time	FEV <sub>1</sub>
001	9.34	1.20	0.00	1.24
001	10.39	1.28	1.05	1.45
001	11.45	1.33	2.11	1.63
001	12.46	1.42	3.12	2.12
001	13.42	1.48	4.08	2.30
001	15.47	1.50	6.13	2.44
001	16.37	1.52	7.03	2.39
002	6.58	1.13	0.00	1.36
002	7.65	1.19	1.06	1.42
002	12.74	1.49	6.15	2.13
002	13.77	1.53	7.19	2.38
002	14.69	1.55	8.11	2.85
002	15.82	1.56	9.23	3.17
002	16.67	1.57	10.08	2.52
002	17.63	1.57	11.04	3.11
003	6.91	1.18	0.00	1.54
003	7.97	1.23	1.06	1.47
003	8.97	1.30	2.05	1.82
003	9.99	1.35	3.08	2.12
003	11.08	1.47	4.16	2.63
003	13.07	1.57	6.16	2.45
003	14.10	1.59	7.19	2.77
003	15.08	1.60	8.17	3.02
003	16.02	1.60	9.10	2.96
007	6.43	1.18	0.00	0.97
007	7.50	1.25	1.06	1.10
007	13.63	1.64	7.19	2.62
007	14.56	1.67	8.12	2.53
007	15.64	1.68	9.21	2.76
007	16.50	1.69	10.06	2.80
007	17.49	1.69	11.06	2.67

Note: Time represents time since entry to study.

When age is used as the metameeter for time, there are two sources of information about the relationship between FEV<sub>1</sub> and age. First, there is “cross-sectional” or between-subject information that arises because children enter the study at different ages. For example, there is information about how FEV<sub>1</sub> changes with age in the baseline (or time = 0) observations. Second, there is “longitudinal” or within-subject information that arises because children are measured repeatedly over time, yielding measurements of FEV<sub>1</sub> at different ages. Because there are two potentially conflicting sources of information about the relationship between FEV<sub>1</sub> and age, it is important to model these data in a way that allows for separate estimation of the “cross-sectional” and “longitudinal” effects of age of FEV<sub>1</sub>. In doing so, it is then possible to test whether there are differences between the cross-sectional and longitudinal effects of age on FEV<sub>1</sub>, and report separate effects where necessary or estimate a combined effect, based on both sources of information, if appropriate. Note that the same issues arise in examining the relationship between FEV<sub>1</sub> and height. A more detailed discussion of the main issues surrounding the analysis of longitudinal designs that provide both longitudinal and cross-sectional sources of information can be found in Chapter 9 (see Sections 9.4 and 9.5).

The Six Cities Study was designed to characterize lung function growth between the ages of 6 and 18. The goal of the following analyses is to determine how changes in lung function (as determined by FEV<sub>1</sub>) over time are related to the age and height of the child. Previous research has indicated that log(FEV<sub>1</sub>) has an approximately linear relationship with age and log(height) in children and adolescents. To distinguish between the cross-sectional and longitudinal effects of age and log(height) on log(FEV<sub>1</sub>), baseline and current values of these covariates were included in the model for the mean. Because these data are inherently unbalanced, accounting for the covariance among the repeated observation on the same child via a random effects structure is very appealing. Here we allow the intercept and slope for age to vary randomly from one child to another. Specifically, we consider the following model for log (FEV<sub>1</sub>):

$$E(Y_{ij}|b_i) = \beta_1 + \beta_2 \text{Age}_{ij} + \beta_3 \log(\text{Ht})_{ij} + \beta_4 \text{Age}_{i1} + \beta_5 \log(\text{Ht})_{i1} + b_{1i} + b_{2i} \text{Age}_{ij},$$

where  $Y_{ij}$  is the log(FEV<sub>1</sub>) for the  $i^{th}$  child at the  $j^{th}$  occasion, and  $\text{Age}_{i1}$  and  $\log(\text{Ht})_{i1}$  are the initial or baseline age and log(height) for the  $i^{th}$  child. In this model,  $\beta_2$  and  $\beta_3$  are the longitudinal effects of

age and  $\log(\text{height})$ , respectively, while  $(\beta_2 + \beta_4)$  and  $(\beta_3 + \beta_5)$  are the corresponding cross-sectional effects. That is,  $\beta_4$  and  $\beta_5$  represent the differences between the longitudinal and cross-sectional effects of age and  $\log(\text{height})$ , respectively. (See Sections 9.4 and 9.5 for a more detailed discussion of, and an alternative representation of, models that allow for estimation of both longitudinal and cross-sectional effects.)

Preliminary analysis of the data revealed a measurement of  $\text{FEV}_1$  that was clearly outlying. This observation was from a girl who had only a baseline measurement available. The observation was removed, and all subsequent analyses are based on the data from 299 girls (with a total of 1993 measurements). The REML estimates of the fixed effects are displayed in [Table 8.2](#). Based on the magnitude of the estimates of  $\beta_4$  and  $\beta_5$ , relative to their standard errors, there is evidence of a significant difference between the longitudinal and cross-sectional effects of age, but not of  $\log(\text{height})$ . From a subject-matter point of view, the magnitudes of the longitudinal and cross-sectional effects of  $\log(\text{height})$  are quite similar (2.24 versus 2.46), whereas the magnitudes of the longitudinal and cross-sectional effects of age are strikingly different (0.024 versus 0.007). That is, the longitudinal and cross-sectional effects of age on changes in  $\text{FEV}_1$  ( $e^{0.024} \approx 1.025$  versus  $e^{0.007} \approx 1.007$ ) are discernibly different. This may be due, in part, to the relatively small amount of variability in ages at baseline (relative to the variability in ages throughout the duration of the study), resulting in the cross-sectional effect of age being poorly estimated from the data at baseline alone; in Sections 9.4 and 9.5 we present an alternative model where estimation of the cross-sectional effect of age is based on measurements at all occasions. From the longitudinal effects of age and  $\log(\text{height})$ , there is clear evidence that changes in  $\log(\text{FEV}_1)$  are related to both age and height.

**Table 8.2** Estimated regression coefficients (fixed effects) and standard errors for the  $\log(\text{FEV}_1)$  data from the Six Cities Study.

Variable	Estimate	SE	Z
Intercept	-0.2883	0.0387	-7.45
Age	0.0235	0.0014	16.86
Log(Height)	2.2372	0.0435	51.39
Initial Age	-0.0165	0.0075	-2.21
Initial Log(Height)	0.2182	0.1455	1.50

Next we consider the interpretation of the fixed effects estimates. The model for the mean, averaged over the distribution of the subject-specific random effects, is given by

$$E(Y_{ij}) = \beta_1 + \beta_2 \text{Age}_{ij} + \beta_3 \log(\text{Ht})_{ij} + \beta_4 \text{Age}_{i1} + \beta_5 \log(\text{Ht})_{i1}.$$

Furthermore this model can be re-expressed in terms of two models, a cross-sectional model and a longitudinal model. The former is given by

$$\begin{aligned} E(Y_{i1}) &= \beta_1 + \beta_2 \text{Age}_{i1} + \beta_3 \log(\text{Ht})_{i1} + \beta_4 \text{Age}_{i1} + \beta_5 \log(\text{Ht})_{i1} \\ &= \beta_1 + (\beta_2 + \beta_4) \text{Age}_{i1} + (\beta_3 + \beta_5) \log(\text{Ht})_{i1}, \end{aligned}$$

while the latter is given by

$$\begin{aligned} E(Y_{ij} - Y_{i1}) &= \beta_1 + \beta_2 \text{Age}_{ij} + \beta_3 \log(\text{Ht})_{ij} + \beta_4 \text{Age}_{i1} + \beta_5 \log(\text{Ht})_{i1} \\ &\quad - \{\beta_1 + \beta_2 \text{Age}_{i1} + \beta_3 \log(\text{Ht})_{i1} + \beta_4 \text{Age}_{i1} + \beta_5 \log(\text{Ht})_{i1}\} \\ &= \beta_2 (\text{Age}_{ij} - \text{Age}_{i1}) + \beta_3 \{\log(\text{Ht})_{ij} - \log(\text{Ht})_{i1}\}. \end{aligned}$$

We note that there are alternative ways to decompose the cross-sectional and longitudinal effects; this topic is explored in greater detail in Chapter 9 (see Sections 9.4 and 9.5).

The longitudinal effect of  $\log(\text{height})$ ,  $\beta_3$ , has interpretation in terms of the changes in mean  $\log(\text{FEV}_1)$  for a single-unit increase in  $\log(\text{height})$ , for any given change in age (e.g., during a two-year interval). Similarly the longitudinal effect of age,  $\beta_2$ , has interpretation in terms of the changes in mean  $\log(\text{FEV}_1)$  for a one-year increase in age, for any given change in  $\log(\text{height})$ . The coefficient for  $\log(\text{height})$ , 2.24, is not directly interpretable because a single-unit change in  $\log(\text{height})$  corresponds to an almost threefold (or  $e^{1.0} \approx 2.7$ ) increase in height. Instead, it is probably more

meaningful to consider the effect of a 10% increase in height. On the logarithmic scale this corresponds to a 0.1 increase in  $\log(\text{height})$ , since  $e^{0.1} \approx 1.1$ . Therefore a 10% increase in height (corresponding to an approximate 0.1 increase in  $\log(\text{height})$ ) is associated with a 0.224 increase in  $\log(\text{FEV}_1)$ . Note that a 0.224 increase in  $\log(\text{FEV}_1)$  corresponds to a 25% increase in  $\text{FEV}_1$  (since  $e^{0.224} = 1.25$ ). On the other hand, the coefficient for age, 0.024, is more directly interpretable. The estimate of the longitudinal effect of age indicates that a single year increase in age is associated with a 0.024 increase in  $\log(\text{FEV}_1)$  or an approximate 2.5% ( $e^{0.024} \approx 1.025$ ) increase in  $\text{FEV}_1$ , for any fixed change in height.

Next consider the estimates<sup>4</sup> of the variances and covariances of the random effects (see [Table 8.3](#)). The marginal covariance among the repeated observations is a function of these variance and covariance parameters (and  $\sigma^2$ ) and the ages of the child when the observations were obtained. The estimated correlations for annual measurements from ages 7 to 18 are displayed in [Table 8.4](#), and these results indicate that there is strong positive correlation among measurements of  $\log(\text{FEV}_1)$  that declines by a modest amount over the 11 years of follow-up. This pattern of correlation reinforces an observation that we made in earlier chapters of the book: the correlation among repeated measurements of many health outcomes rarely decays to zero, even when they are separated by many years.

**Table 8.3** Estimated covariance of the random effects and standard errors (x 100) for the  $\log(\text{FEV}_1)$  data from the Six Cities Study.

Parameter	Estimate	SE
$\text{Var}(b_{1i}) = g_{11}$	1.2207	0.1924
$\text{Cov}(b_{1i}, b_{2i}) = g_{12}$	-0.0435	0.0122
$\text{Var}(b_{2i}) = g_{22}$	0.0050	0.0010
$\text{Var}(\epsilon_i) = \sigma^2$	0.3629	0.0133

**Table 8.4** Estimated marginal correlations among repeated measures of  $\log(\text{FEV}_1)$  between the ages of 7 and 18.

Age (years)												
7	8	9	10	11	12	13	14	15	16	17	18	
1.00	0.70	0.69	0.68	0.67	0.66	0.64	0.62	0.60	0.58	0.56	0.54	
0.70	1.00	0.70	0.69	0.68	0.67	0.66	0.65	0.63	0.61	0.60	0.58	
0.69	0.70	1.00	0.70	0.70	0.69	0.68	0.67	0.66	0.64	0.63	0.61	
0.68	0.69	0.70	1.00	0.70	0.70	0.70	0.69	0.68	0.67	0.66	0.64	
0.67	0.68	0.70	0.70	1.00	0.71	0.71	0.70	0.70	0.69	0.68	0.67	
0.66	0.67	0.69	0.70	0.71	1.00	0.72	0.72	0.71	0.71	0.70	0.70	
0.64	0.66	0.68	0.70	0.71	0.72	1.00	0.73	0.73	0.73	0.72	0.72	
0.62	0.65	0.67	0.69	0.70	0.72	0.73	1.00	0.74	0.74	0.74	0.74	
0.60	0.63	0.66	0.68	0.70	0.71	0.73	0.74	1.00	0.75	0.75	0.75	
0.58	0.61	0.64	0.67	0.69	0.71	0.73	0.74	0.75	1.00	0.76	0.76	
0.56	0.60	0.63	0.66	0.68	0.70	0.72	0.74	0.75	0.76	1.00	0.77	
0.54	0.58	0.61	0.64	0.67	0.70	0.72	0.74	0.75	0.76	0.77	1.00	

Finally, we note that the correlation among repeated measurements has been accounted for by the introduction of random intercepts and slopes for age. Alternatively, we could have considered a random effects model with randomly varying slopes for  $\log(\text{height})$ . By assuming that the slope for  $\log(\text{height})$  varies randomly for individuals, we can also induce covariance among the repeated observations but with correlations that are a function not of age, but of the height of the child. For illustrative purposes we considered the following model:

$$\begin{aligned} E(Y_{ij}|b_i) &= \beta_1 + \beta_2 \text{Age}_{ij} + \beta_3 \log(\text{Ht})_{ij} + \beta_4 \text{Age}_{i1} + \beta_5 \log(\text{Ht})_{i1} \\ &\quad + b_{1i} + b_{2i} \log(\text{Ht})_{ij}. \end{aligned}$$

The REML estimates of the fixed effects are displayed in [Table 8.5](#) and are qualitatively very similar to those presented in [Table 8.2](#). The reader might then ask which model is more appropriate for the

data at hand, a model with randomly varying slopes for age or randomly varying slopes for log(height)? Fortunately, since both models have the same number of covariance parameters, we can make that judgement based on a direct comparison of their maximized REML log-likelihoods. For the model with randomly varying slopes for age the maximized REML log-likelihood is 2283.9, while for the model with randomly varying slopes for log(height) the maximized REML log-likelihood is 2294.7. The comparison of the maximized log-likelihoods indicates that the model with randomly varying slopes for log(height) is to be preferred. For illustrative purposes we also considered a random effects model with randomly varying slopes for both age and log(height). By assuming that the slopes for age and log(height) vary randomly this would induce covariances among the repeated observations that are functions of both the age and height of the child. For the latter model the maximized REML log-likelihood is 2294.9 and does not lead to a discernible improvement in fit over the model with randomly varying slopes for log(height) only.

**Table 8.5** Estimated regression coefficients (fixed effects) and standard errors for the log(FEV<sub>1</sub>) data from the Six Cities Study.

Variable	Estimate	SE	Z
Intercept	-0.2846	0.0390	-7.30
Age	0.0233	0.0012	18.65
Log(Height)	2.2523	0.0461	48.82
Initial Age	-0.0163	0.0074	-2.19
Initial Log(Height)	0.1808	0.1455	1.24

Formally, the likelihood ratio test statistic,  $G^2 = 0.4$ , can be compared to the critical values in the third row ( $q = 2$ ) of [Table C.1](#) in Appendix C.

# Study of Influence of Menarche on Changes in Body Fat Accretion

The second illustration uses longitudinal data from a prospective study on body fat accretion in a cohort of 162 girls from the MIT Growth and Development Study (Bandini et al., 2002; Phillips et al., 2003). The data represent a subset of the study materials and should not be used to draw substantive conclusions.

It is known that increases in body fatness in girls begin just before or around menarche. Although it has been presumed that the increase in body fatness levels off approximately four years after menarche, these changes in body fat accretion had not been studied in population-based samples. Naumova et al. (2001) examined changes in body fat before and after menarche. At the start of the study, all of the girls were pre-menarcheal and non-obese, as determined by a triceps skinfold thickness less than the 85th percentile. All girls were followed over time according to a schedule of annual measurements until four years after menarche. The final measurement was scheduled on the fourth anniversary of their reported date of menarche.

At each examination, a measure of body fatness was obtained based on bioelectric impedance analysis. Percent body fat (%BF) was derived from three basic measurements of body weight (Wt. in kg), height (Ht. in cm), and bioelectric impedance resistance (R). Percent body fat is calculated using the equation:

$$\%BF = \left( 1 - \frac{TBW}{0.73} Wt \right) \times 100\%,$$

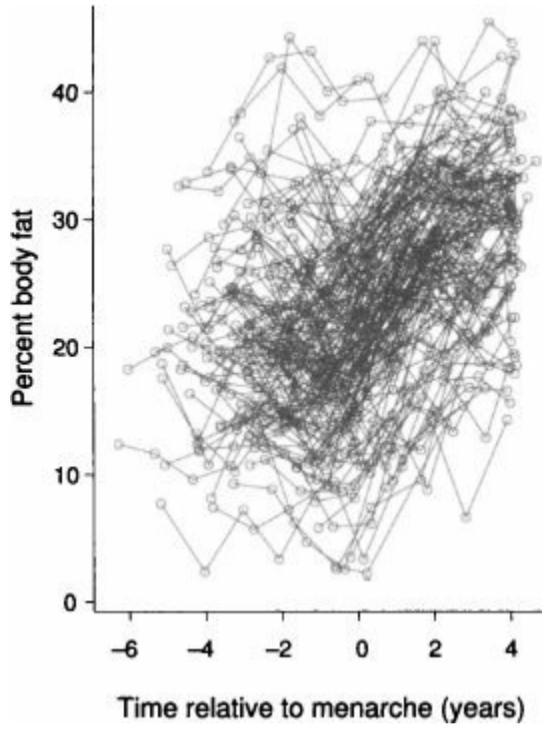
where total body water,  $TBW = (0.7Ht^2/R) - 0.32$ .

In this section we present an analysis of the changes in percent body fat before and after menarche. For the purposes of these analyses, “time” is coded as time since menarche and can be positive or negative. Although the measurement protocol is the same for all girls and the study design is balanced if the timing of measurement is defined as the time since the baseline measurement, it is inherently unbalanced when the timing of measurements is defined as the time since a girl experienced menarche.

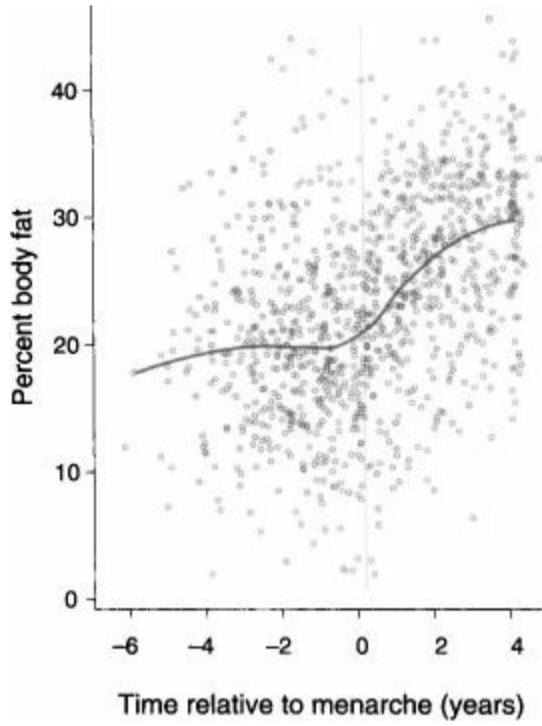
In this data set there are a total of 1049 percent body fat measurements, with an average of 6.5 measurements per subject. The numbers of measurements per subject pre- and post-menarche are approximately equal, with 497 measurements for the pre-menarcheal period (producing an average of 3.1 measurements per subject) and 552 measurements for the post-menarcheal period (producing an average of 3.5 measurements per subject). In this sample the average age at menarche was 12.8 years.

[Figure 8.5](#) shows a time plot of the individual response profiles (where time is relative to the individual age at menarche). This graph reveals some information about the greater variability of measurement times before menarche. However, it is difficult to discern whether the changes in percent body fat in the pre-menarcheal period are similar to the changes in the post-menarcheal period. In [Figure 8.6](#) the trend in the mean response is assessed using a *lowess* smoothed curve. Recall that *lowess* is a nonparametric, robust regression method that traces the salient features of the mean response as a function of time while making only minimal assumptions about the form of the relationship. The *lowess* curve reveals that the mean response remains relatively flat during the pre-menarcheal period and then rises sharply during the post-menarcheal period.

[\*\*Fig. 8.5\*\*](#) Time plot of percent body fat against time, relative to age of menarche (in years).



**Fig. 8.6** Time plot of percent body fat against time, relative to age of menarche (in years), with lowess smoothed curve.



In the following analysis we consider the hypothesis that percent body fat accretion increases linearly with age, but with different slopes before and after menarche. Specifically, we assume that each girl has a piecewise linear spline growth curve with a knot at the time of menarche. That is, each girl's growth curve can be described with an intercept and two slopes, one slope for the changes in response before menarche, another slope for the changes in response after menarche. Note that unlike the piecewise linear splines considered in Section 6.3, the knot is not the same age for all subjects.

Let  $t_{ij}$  denote the time of the  $j^{th}$  measurement on the  $i^{th}$  subject before or after menarche (i.e.,  $t_{ij} = 0$  at menarche). We fit the following linear mixed effects model:

$$E(Y_{ij}|b_i) = \beta_1 + \beta_2 t_{ij} + \beta_3 (t_{ij})_+ + b_{1i} + b_{2i} t_{ij} + b_{3i} (t_{ij})_+,$$

where  $(t_{ij})_+ = t_{ij}$  if  $t_{ij} > 0$  and  $(t_{ij})_+ = 0$  if  $t_{ij} \leq 0$ . In this model,  $(\beta_1 + b_{1i})$  is the intercept for the  $i^{th}$  subject and has interpretation as the true percent body fat at menarche (when  $t_{ij} = 0$ ). Of note, the actual percent body fat at menarche is not observed and cannot be directly estimated from the data at hand. As a result we use the term "true" percent body fat at menarche to remind the reader that this is a parameter in the model. Similarly  $(\beta_2 + b_{2i})$  is the  $i^{th}$  subject's slope, or rate of change in percent body fat during the pre-menarcheal period. Finally, the  $i^{th}$  subject's slope during the post-menarcheal period is given by  $\{(\beta_2 + \beta_3) + (b_{2i} + b_{3i})\}$ . Since the overall goal of the analysis is to assess whether the population slopes for fat accretion differ before and after menarche, this can be translated into the

null hypothesis,  $H_0: \beta_3 = 0$ .

The REML estimates of the fixed effects and the variance components are displayed in [Tables 8.6](#) and [8.7](#), respectively. Based on the magnitude of the estimate of  $\beta_3$ , relative to its standard error, there is a significant difference between the slopes before and after menarche. The estimate of the population mean pre-menarcheal slope is 0.42, which is statistically significant at the 0.05 level. This estimated slope is rather shallow and indicates that the annual rate of body fat accretion is less than 0.5%. Note that the estimated variance of  $b_{2i}$  is 1.63, indicating that there is substantial variability from girl to girl in rates of fat accretion and that many girls are losing body fat while others are gaining body fat during the pre-menarcheal period. For example, approximately 95% of girls have changes in percent body fat between  $-2.09\%$  and  $2.92\%$  (i.e.,  $0.42 \pm 1.96 \times \sqrt{1.63}$ ). The estimate of the population mean post-menarcheal slope is 2.46 (with SE = 0.12), which is statistically significant at the 0.05 level. This indicates that the annual rate of body fat accretion is approximately 2.5%, almost six times higher than the corresponding rate in the pre-menarcheal period. The estimated variance of the individual slopes during the post-menarcheal period,  $\text{Var}(b_{2i} + b_{3i})$ , is 0.88 (or  $[1.63 + 2.75 - 2 \times 1.75]$ ), indicating that there is less variability in the slopes after menarche. For example, approximately 95% of girls have changes in percent body fat between  $0.62\%$  and  $4.30\%$  (i.e.,  $2.46 \pm 1.96 \times \sqrt{0.88}$ ). In other words, more than 95% of girls are expected to have increases in body fat during the post-menarcheal period, while substantially fewer (approximately 63%) are expected to have increases in body fat during the pre-menarcheal period.

**Table 8.6** Estimated regression coefficients (fixed effects) and standard errors for the percent body fat data.

Variable	Estimate	SE	Z
Intercept	21.3614	0.5646	37.84
Time	0.4171	0.1572	2.65
(Time) <sub>+</sub>	2.0471	0.2280	8.98

**Table 8.7** Estimated covariance of the random effects and standard errors for the percent body fat data.

Parameter	Estimate	SE
$\text{Var}(b_{1i}) = g_{11}$	45.9413	5.7393
$\text{Var}(b_{2i}) = g_{22}$	1.6311	0.4331
$\text{Var}(b_{3i}) = g_{33}$	2.7497	0.9635
$\text{Cov}(b_{1i}, b_{2i}) = g_{12}$	2.5263	1.2185
$\text{Cov}(b_{1i}, b_{3i}) = g_{13}$	-6.1096	1.8730
$\text{Cov}(b_{2i}, b_{3i}) = g_{23}$	-1.7505	0.5980
$\text{Var}(\epsilon_i) = \sigma^2$	9.4732	0.5443

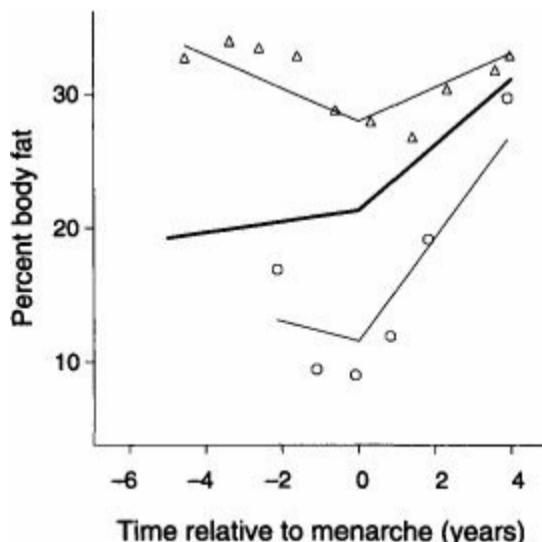
The estimated marginal correlations among annual measurements of percent body fat, based on the estimated covariances among the random effects, are displayed in [Table 8.8](#). These results indicate that there is strong positive correlation (approximately 0.8) between annual measurements of percent body fat. Although the strength of the correlation declines over time, it does not decay to zero even when measurements are taken eight years apart. In general, the variability of percent body fat is greater in the pre-menarcheal period.

**Table 8.8** Estimated marginal correlations (on the off-diagonals) among repeated measures of percent body fat between 4 years pre- and post-menarche, with estimated variances along the main diagonal.

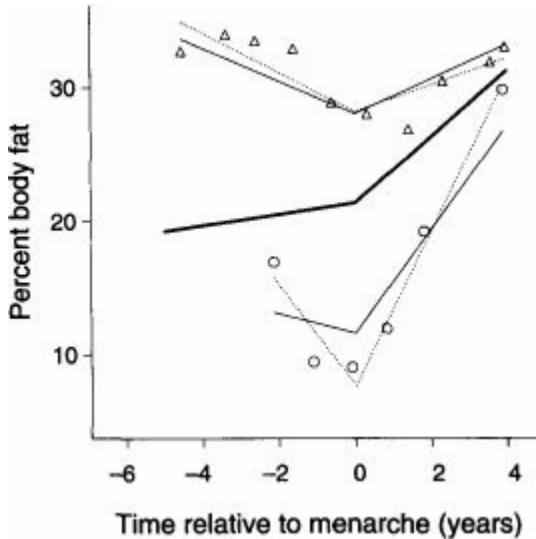
Time (relative to menarche)									
-4	-3	-2	-1	0	1	2	3	4	
61.3	0.82	0.78	0.71	0.61	0.60	0.57	0.52	0.47	
0.82	54.9	0.81	0.76	0.70	0.68	0.64	0.60	0.54	
0.78	0.81	51.8	0.80	0.76	0.74	0.71	0.66	0.60	
0.71	0.76	0.80	52.0	0.81	0.79	0.76	0.71	0.64	
0.61	0.70	0.76	0.81	55.4	0.81	0.78	0.73	0.66	
0.60	0.68	0.74	0.79	0.81	49.1	0.79	0.76	0.70	
0.57	0.64	0.71	0.76	0.78	0.79	44.6	0.77	0.74	
0.52	0.60	0.66	0.71	0.73	0.76	0.77	41.8	0.76	
0.47	0.54	0.60	0.64	0.66	0.70	0.74	0.76	40.8	

The mixed effects model can be used to obtain estimates of each girl's growth trajectory over time. [Figure 8.7](#) displays the estimated population mean growth curve and the predicted (empirical BLUP) growth curves for two girls, based on the fixed and random effects estimates reported in [Tables 8.6](#) and [8.7](#). Note that the two girls selected for display in [Figure 8.7](#) differ in the number of measurements that were obtained (with 6 and 10 measurements, respectively). A noticeable feature of the predicted growth curves is that there is more shrinkage toward the population mean curve when fewer data points are available. That is, the predicted growth curve for the girl with only 6 data points is pulled closer to the population mean curve (or further away from her own data points) while the predicted growth curve for the girl with 10 observation follows her data more closely. This feature becomes more apparent when the empirical BLUPs are compared to the ordinary least squares (OLS) estimates based only on the longitudinal observations from each girl (see [Figure 8.8](#)). Examination of [Figure 8.8](#) reveals that the empirical BLUP for the girl with 10 observations is largely based on her longitudinal observations. On the other hand, the empirical BLUP for the girl with 6 observations "borrows strength" from the population mean curve. This is a characteristic feature of the empirical BLUPs that was noted in Sections 8.6 and 8.7. When there is less information available for estimating an individual's growth curve, there is a greater "borrowing of strength" from the data obtained on all girls in the study.

**Fig. 8.7** Population average curve (thicker solid line) and empirical BLUPs for two randomly selected girls.



**Fig. 8.8** Population average curve, empirical BLUPs, and OLS predictions for two randomly selected girls.



Finally, we can use these data to illustrate a hybrid random effects and covariance pattern model by fitting the following model to the percent body fat:

$$Y_{ij} = \beta_1 + \beta_2 t_{ij} + \beta_3(t_{ij})_+ + b_{1i} + U_i(t_{ij}) + \epsilon_{ij},$$

where the  $U_i(t_{ij})$  are assumed to have a normal distribution, with zero mean, variance  $\sigma^2_u$ , and correlation

$$\text{Corr}\{U_i(t_{ij}), U_i(t_{ik})\} = \rho(|t_{ij} - t_{ik}|).$$

The  $U_i(t_{ij})$  induce serial correlation among the responses, such that the correlation becomes weaker as the time separation increases. Two popular choices for  $\rho(|t_{ij} - t_{ik}|)$  are the exponential correlation function,

$$\rho(|t_{ij} - t_{ik}|) = e^{-\alpha|t_{ij} - t_{ik}|},$$

and the Gaussian correlation function,

$$\rho(|t_{ij} - t_{ik}|) = e^{-\alpha|t_{ij} - t_{ik}|^2},$$

for some  $\alpha > 0$ . Finally, the  $\epsilon_{ij}$  are the usual sampling or measurement errors and these are assumed to be independent, with mean zero and variance  $\sigma^2$ .

**Table 8.9** Comparison of the maximized (REML) log-likelihoods and AIC for the mixed effects model and the hybrid models with exponential and Gaussian serial correlation.

Model	-2 (REML) Log-Likelihood	AIC
Mixed Effects	6062.4	6076.4
Hybrid: Exponential Serial Correlation	5999.9	6007.9
Hybrid: Gaussian Serial Correlation	5991.2	5999.2

**Table 8.10** Estimated regression coefficients (fixed effects) and standard errors for the hybrid model with Gaussian serial correlation.

Variable	Estimate	SE	Z
Intercept	21.2918	0.5400	39.43
Time	0.2168	0.1439	1.51
(Time) <sub>+</sub>	2.1655	0.2331	9.29

We considered the goodness of fit of the hybrid model when the serial correlation function is exponential and Gaussian. [Table 8.9](#) displays the maximized (REML) log-likelihood and AIC for the hybrid models with exponential and Gaussian serial correlation; the maximized (REML) log-likelihood and AIC for the mixed effects model considered previously are also displayed. These results indicate that the hybrid model with Gaussian serial correlation fits the data best, since it has the largest maximized log-likelihood (when compared to the hybrid model with exponential serial correlation) and the smallest AIC (when compared to the mixed effects model).

The REML estimates of the fixed effects from the hybrid model with Gaussian serial correlation are displayed in [Table 8.10](#). The estimates of  $\beta$  are similar to those reported in [Table 8.6](#). In particular, the estimate of  $\beta_3$  is very similar, and when compared to its standard error, there is a significant difference between the slopes before and after menarche. On the other hand, the estimate

of the population mean pre-menarcheal slope is 0.22, and is no longer statistically significant at the 0.05 level. Overall, the substantive conclusions are very similar in the two sets of analyses: there is at most a very weak pre-menarcheal slope, indicating that the annual rate of body rate accretion is very modest (0.2–0.4%), while the annual rate of fat accretion during the post-menarcheal period is discernibly greater (approximately 2.4–2.5%) than the corresponding rate in the pre-menarcheal period. Of note, an attempt to fit an extended mixed effects model (with randomly varying intercepts and pre- and post-menarcheal slopes) by incorporation of a Gaussian serial correlation component failed to converge. This lack of convergence was taken as an indication that the observed data simply do not support the need for both randomly varying slopes and serially correlated residuals. As was mentioned in Section 8.2, there can be identifiability problems with the hybrid model unless the random effects structure is kept very simple (e.g., random intercepts only). That is, there may be insufficient information in the data at hand to support separate estimation of randomly varying slopes, serially correlated residuals, and measurement errors.

# Randomized Study of Dual or Triple Combinations of HIV-1 Reverse Transcriptase Inhibitors

The final illustration uses data from a randomized, double-blind, study of AIDS patients with advanced immune suppression (CD4 counts of  $\leq 50$  cells/mm $^3$ ) (Henry et al., 1998). Patients in AIDS Clinical Trial Group (ACTG) Study 193A were randomized to dual or triple combinations of HIV-1 reverse transcriptase inhibitors. Specifically, patients were randomized to one of four daily regimens containing 600 mg of zidovudine: zidovudine alternating monthly with 400 mg didanosine, zidovudine plus 2.25 mg of zalcitabine, zidovudine plus 400 mg of didanosine, or zidovudine plus 400 mg of didanosine plus 400 mg of nevirapine (triple therapy). For the analyses presented here, we focus on the comparison of the first three treatment regimens (dual therapy) with the fourth (triple therapy).

Measurements of CD4 counts were scheduled to be collected at baseline and at 8-week intervals during follow-up. However, the CD4 count data are unbalanced due to mistimed measurements and missing data that resulted from skipped visits and dropout (see data for four randomly selected subjects presented in [Table 8.11](#)). The number of measurements of CD4 counts during the first 40 weeks of follow-up varied from 1 to 9, with a median of 4. The goal of our analyses is to compare the dual and triple therapy groups in terms of short-term changes in CD4 counts from baseline to week 40 (approximately 10 months of follow-up). The analyses are based on log transformed CD4 counts,  $\log(\text{CD4 count}+1)$ , available on 1309 patients.

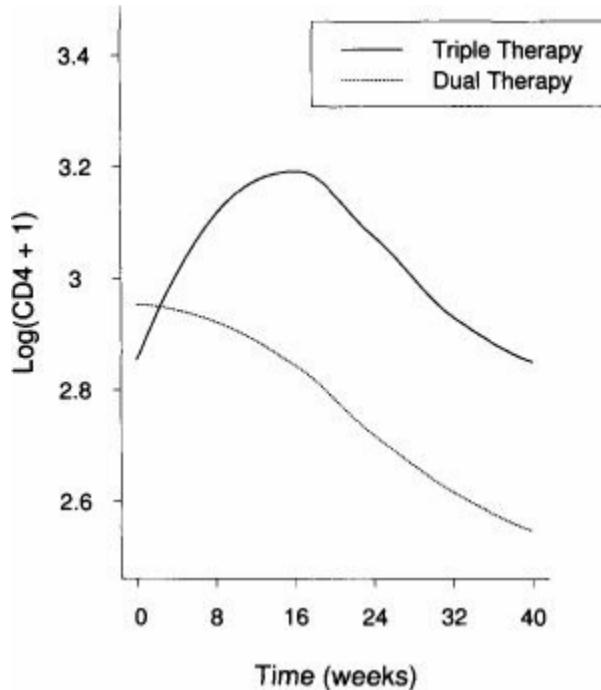
**Table 8.11** Data on log CD4 counts for four randomly selected subjects from ACTG study 193A.

Subject ID	Group	Time	$\log(\text{CD4} + 1)$
56	0	0.0	1.7047
56	0	8.1	1.7918
56	0	16.1	0.6932
56	0	25.4	1.0986
56	0	33.4	0.6932
56	0	39.1	0.6932
544	1	0.0	3.3844
544	1	7.6	3.2189
544	1	15.9	2.1972
544	1	31.9	1.6094
736	0	0.0	3.7495
736	0	8.9	3.4965
736	0	18.9	3.1780
736	0	30.9	2.7726
986	1	0.0	4.4659
986	1	17.4	3.3322
986	1	30.9	3.5553
986	1	39.6	3.3673

*Note:* Group = 1 if randomized to triple therapy, Group = 0 if randomized to dual therapy.

In [Figure 8.9](#) the trend in the mean response in the dual and triple therapy groups is assessed using *lowess* smoothed curves. The curves reveal a modest decline in the mean response during the first 16 weeks for the dual therapy group, followed by a steeper decline from week 16 to week 40. In contrast, for the triple therapy group, the mean response increases during the first 16 weeks and declines thereafter. The rate of decline from week 16 to week 40 appears to be similar for the two groups. A note of caution: because there is a substantial amount of missing data, the plot of the mean response over time can be potentially misleading unless the data are missing completely at random (MCAR). When data are missing, and especially when there is attrition of subjects whose responses are different from those who remain in the study, a plot of the mean response over time can be deceptive; the observed changes in the mean response may reflect the pattern of missingness or the attrition, and not within-individual change. (See Chapters 17 and 18 for a more detailed discussion of this issue.)

**Fig. 8.9** Lowess smoothed curves of  $\log(\text{CD4} + 1)$  against time (in weeks), for subjects in the dual and triple therapy groups in ACTG study 193A.



Next we consider a model for the mean response that allows the rates of change before and after week 16 to differ within and between groups. Specifically, we assume that each patient has a piecewise linear spline with a knot at week 16. That is, each patient's response trajectory can be described with an intercept and two slopes—one slope for the changes in response before week 16, another slope for the changes in response after week 16. The average slopes for changes in response before and after week 16 are allowed to vary by group. Because this is a randomized study, the mean response at baseline is assumed to be the same in the two groups.

Letting  $t_{ij}$  denote the time since baseline for the  $j^{th}$  measurement on the  $i^{th}$  subject (with  $t_{ij} = 0$  at baseline), we consider the following linear mixed effects model:

$$E(Y_{ij}|b_i) = \beta_1 + \beta_2 t_{ij} + \beta_3(t_{ij} - 16)_+ + \beta_4 \text{Group}_i \times t_{ij} + \beta_5 \text{Group}_i \times (t_{ij} - 16)_+$$

$$+ b_{1i} + b_{2i} t_{ij} + b_{3i} (t_{ij} - 16)_+,$$

where  $\text{Group}_i = 1$  if the  $i^{th}$  subject is randomized to triple therapy, and  $\text{Group}_i = 0$  otherwise;  $(t_{ij} - 16)_+ = t_{ij} - 16$  if  $t_{ij} > 16$  and  $(t_{ij} - 16)_+ = 0$  if  $t_{ij} \leq 16$ . In this model,  $(\beta_1 + b_{1i})$  is the intercept for the  $i^{th}$  subject and has interpretation as the true log CD4 count at baseline (when  $t_{ij} = 0$ ). Similarly  $(\beta_2 + b_{2i})$  is the  $i^{th}$  subject's slope, or rate of change in log CD4 counts from baseline to week 16, if randomized to dual therapy;  $(\beta_2 + \beta_4 + b_{2i})$  is the  $i^{th}$  subject's slope if randomized to triple therapy. Finally, the  $i^{th}$  subject's slope from week 16 to week 40 is given by  $\{(\beta_2 + \beta_3) + (b_{2i} + b_{3i})\}$  if randomized to dual therapy and  $\{(\beta_2 + \beta_3 + \beta_4 + \beta_5) + (b_{2i} + b_{3i})\}$  if randomized to triple therapy. The null hypothesis of no treatment group differences in the changes in log CD4 counts can be expressed as  $H_0: \beta_4 = \beta_5 = 0$ .

The REML estimates of the fixed effects are displayed in [Table 8.12](#). A test of  $H_0: \beta_4 = \beta_5 = 0$  yields a Wald test,  $W^2 = 59.21$ , with 2 degrees of freedom ( $p < 0.0001$ ); the corresponding likelihood ratio test yields  $G^2 = 57.99$ , with 2 degrees of freedom ( $p < 0.0001$ ). Based on the magnitude of the estimate of  $\beta_4$ , relative to its standard error, there is a significant group difference in the rates of change from baseline to week 16. In the dual therapy group, there is a significant decrease in the mean of the log CD4 counts from baseline to week 16. The estimated change during the first 16 weeks is  $-0.12$ , or  $16 \times -0.0073$ . On the untransformed scale, this corresponds to an approximate 10% decrease in CD4 counts (since  $e^{-0.12} = 0.89$ ). In contrast, in the triple therapy group, there is a significant increase in the mean response. The estimated change during the first 16 weeks in the triple therapy group is  $0.31$ , or  $16 \times (-0.0073 + 0.0269)$ ; the estimated slope for the triple therapy group,  $0.0196$ , has a standard error of  $0.0033$ . On the untransformed scale, this corresponds to an approximate 35% increase in CD4 counts (since  $e^{0.31} = 1.36$ ).

**Table 8.12** Estimated regression coefficients (fixed effects) and standard errors for the log CD4 counts.

Variable	Estimate	SE	Z
Intercept	2.9415	0.0256	114.81
$t_{ij}$	-0.0073	0.0020	-3.70
$(t_{ij} - 16)_+$	-0.0120	0.0032	-3.79
Group $\times t_{ij}$	0.0269	0.0039	6.98
Group $\times (t_{ij} - 16)_+$	-0.0277	0.0062	-4.47

The lowess curves in [Figure 8.9](#) suggest that the rate of decline from week 16 to week 40 is similar for the two groups. The null hypothesis of no treatment group differences in the rates of change in log CD4 counts from week 16 to week 40 can be expressed as  $H_0: \beta_4 + \beta_5 = 0$  (or  $H_0: \beta_4 = -\beta_5$ ). The estimates of  $\beta_4$  and  $\beta_5$  in [Table 8.12](#) appear to support the null hypothesis since they are of similar magnitude but opposite sign. A test of the null hypothesis,  $H_0: \beta_4 + \beta_5 = 0$ , yields a Wald test,  $W^2 = 0.07$ , with 1 degree of freedom ( $p > 0.75$ ); the corresponding likelihood ratio test yields  $G^2 = 0.07$ , with 1 degree of freedom ( $p > 0.75$ ).

The estimated variances of the random effects in [Table 8.13](#) indicate that there is substantial variability from patient to patient in baseline CD4 counts and the rates of change in CD4 counts. For example, although many patients randomized to triple therapy show increases in CD4 counts during the first 16 weeks, some patients have declining CD4 counts. Specifically, approximately 95% of patients randomized to triple therapy are expected to have changes in log CD4 counts from baseline to week 16 between  $-0.64$  and  $1.27$  (or  $16 \times [0.0196 \pm 1.96 \times \sqrt{0.000923}]$ ). That is, approximately 26% of patients are expected to have decreases in CD4 counts during the first 16 weeks of triple therapy. There is also a substantial component of variability due to measurement error.

**Table 8.13** Estimated covariance (x 1000) of the random effects and standard errors for the log CD4 counts.

Parameter	Estimate	SE
$\text{Var}(b_{1i}) = g_{11}$	585.742	34.754
$\text{Var}(b_{2i}) = g_{22}$	0.923	0.160
$\text{Var}(b_{3i}) = g_{33}$	1.240	0.395
$\text{Cov}(b_{1i}, b_{2i}) = g_{12}$	7.254	1.805
$\text{Cov}(b_{1i}, b_{3i}) = g_{13}$	-12.348	2.730
$\text{Cov}(b_{2i}, b_{3i}) = g_{23}$	-0.919	0.236
$\text{Var}(\epsilon_i) = \sigma^2$	306.163	10.074

In a clinical trial it is often of interest to predict the direction and magnitude of the treatment effect for patients with specific covariate values. In the physician–patient context, for example, these predictions can be used to identify those patients who do not respond well to their assigned therapy. When there is interest in subject-specific predictions, we must consider the relative magnitudes of the between-subject and within-subject variability. When the within-subject or measurement error variability is relatively large, the observed response profile for a subject is unreliable and a better prediction can be obtained by “borrowing strength” from the data on all of the subjects. Next we consider the prediction of patients’ response trajectories from the following linear mixed effects model which also includes gender and baseline age:

$$\begin{aligned}
 E(Y_{ij}|b_i) &= \beta_1 + \beta_2 t_{ij} + \beta_3 (t_{ij} - 16)_+ + \beta_4 \text{Group}_i \times t_{ij} - \beta_4 \text{Group}_i \times (t_{ij} - 16)_+ \\
 &\quad + \beta_5 \text{Age}_i + \beta_6 \text{Gender}_i + b_{1i} + b_{2i} t_{ij} + b_{3i} (t_{ij} - 16)_+ \\
 &= \beta_1 + \beta_2 t_{ij} + \beta_3 (t_{ij} - 16)_+ + \beta_4 \text{Group}_i \times \{t_{ij} - (t_{ij} - 16)_+\} \\
 &\quad + \beta_5 \text{Age}_i + \beta_6 \text{Gender}_i + b_{1i} + b_{2i} t_{ij} + b_{3i} (t_{ij} - 16)_+,
 \end{aligned}$$

where  $\text{Age}_i$  is the baseline age (in years) of the patient,  $\text{Gender}_i = 1$  if the  $i^{th}$  patient is male, and

$\text{Gender}_i = 0$  otherwise. In this model the mean rate of change from baseline to week 16 can differ in the two groups (with slopes of  $\beta_2$  and  $\beta_2 + \beta_4$  respectively), but the mean rate of change from week 16 to week 40 is assumed to be the same (with slope of  $\beta_2 + \beta_3$ ).

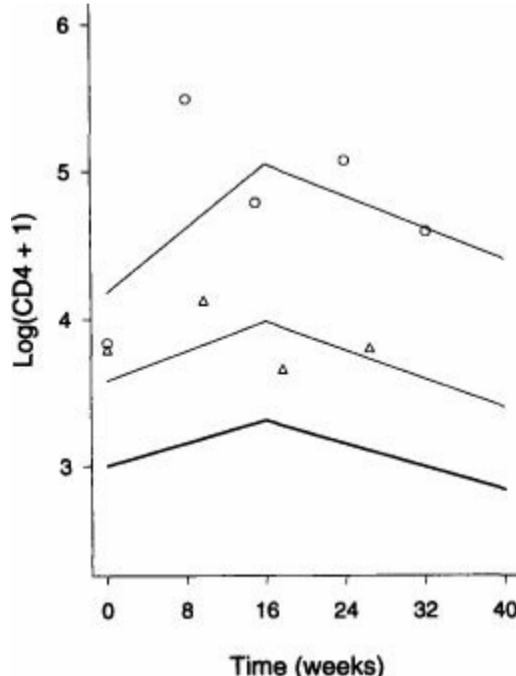
The REML estimates of the fixed effects are displayed in [Table 8.14](#) and the substantive conclusions about the treatment group comparisons are similar to those obtained from [Table 8.12](#). Adjusting for gender and age at baseline, there is a 10% decrease in CD4 counts (since  $e^{16x-0.0072} = e^{-0.12} = 0.89$ ) in the dual therapy group. In contrast, in the triple therapy group, there is a 35% increase in CD4 counts (since  $e^{16x(-0.0072+0.0263)} = e^{0.31} = 1.36$ ). In both treatment groups, there is a significant decline in the mean response from week 16 to week 40, corresponding to an approximate 40% decrease in CD4 counts (since  $e^{24x(-0.0072-0.0124)} = e^{-0.47} = 0.63$ ).

**Table 8.14** Estimated regression coefficients (fixed effects) and standard errors for the revised model for the log CD4 counts.

Variable	Estimate	SE	Z
Intercept	2.6457	0.1280	20.67
$t_{ij}$	-0.0072	0.0019	-3.71
$(t_{ij} - 16)_+$	-0.0124	0.0029	-4.33
Group $\times \{t_{ij} - (t_{ij} - 16)_+\}$	0.0263	0.0034	7.68
Age	0.0100	0.0030	3.31
Gender	-0.0927	0.0754	-1.23

The inclusion of the random effects in the model allows each patient's response trajectory to be described with an intercept and two slopes, one slope for the changes in response before week 16, another slope for the changes in response after week 16. Based on the REML estimates of the fixed effects and variance components, the predicted (or BLUP) trajectory for each patient can be obtained. [Figure 8.10](#) displays the estimated population mean curve and the predicted curves for two male patients, aged 45, and with similar baseline CD4 counts, who were randomized to triple therapy.

**Fig. 8.10** Population average curve (thicker solid line) and empirical BLUPs for two male patients, aged 45, with similar baseline CD4 counts, and treated with triple therapy.



In general, the empirical BLUPs, or the predictions of summary measures (e.g., predicted area under the curve for a patient), can be used to identify those patients who have or have not responded well to their assigned therapy. In the physician–patient context, these predictions may be far more relevant than knowledge of the population mean curve. The appealing feature of the linear mixed effects model analysis is that it allows inferences about both the population trends and individual-specific trajectories.

# 8.9 COMPUTING: FITTING LINEAR MIXED EFFECTS MODELS USING PROC MIXED IN SAS

To fit linear mixed effects models we need to make use of the RANDOM statement in PROC MIXED. The RANDOM statement is used to define all effects that are considered to be random. Specifically, the RANDOM statement is used to define the covariates in the design matrix,  $Z_i$ , for the random effects,  $b_i$ . Ordinarily these will be a subset of the covariates included on the MODEL statement. (Recall that the covariates in the design matrix,  $X_i$ , for the fixed effects appear in the MODEL statement.) While the MODEL statement is used to define the design matrix for the fixed effects and the RANDOM statement is used to define the design matrix for the random effects, note that an intercept is included by default in the former but not the latter. That is, unlike the MODEL statement, PROC MIXED does not include an intercept in the RANDOM statement by default. However, you can specify INTERCEPT (or INT) as a random effect on the RANDOM statement. The RANDOM statement is also used to specify the structure of the covariance matrix for the random effects, G. The structure of G is specified using the TYPE=option. The random effects can be assumed to be correlated (TYPE=UN) or uncorrelated (TYPE=VC); ordinarily, covariance pattern models are not used to account for the covariance among the random effects. For reasons discussed in Section 8.2, we recommend using an unstructured covariance matrix (TYPE=UN) for G. To ensure that the unstructured covariance matrix for the random effects is constrained to be positive-definite, the TYPE=FAO(q) option can be used (where  $q$  is the number of random effects). The latter option can be useful when the TYPE=UN option yields an estimated G matrix that is not positive-definite.

For example, to fit a model with randomly varying intercepts and slopes to data from two or more groups measured repeatedly over time, we can use the illustrative SAS commands given in [Table 8.15](#). Note that the SUBJECT option on the RANDOM statement is used in the same manner as on the REPEATED statement and denotes a variable that distinguishes clusters of correlated responses. By including a variable in the SUBJECT option (e.g., a subject identifier), pairs of observations with distinct values of that variable are regarded as independent. Pairs of observations with the same values of that variable share common values of the random effects.

**Table 8.15** Illustrative commands for a linear mixed effects model, with randomly varying intercepts and slopes, using PROC MIXED in SAS.

---

```
PROC MIXED;
  CLASS id group;
  MODEL y=group time group*time/S CHISQ;
  RANDOM INTERCEPT time/TYPE=UN SUBJECT=id G V;
```

---

Various options can be included on the RANDOM statement. The option G requests that the estimates of the variances and covariances of the random effects be displayed. The option GCORR requests that the estimates of the correlations among the random effects be displayed. The option V requests that the estimates of the marginal covariance matrix, averaged over the distribution of the random effects, be displayed for the first subject. That is, the option V produces estimates of  $\Sigma_i = \text{Cov}(Y_i) = Z_i G Z_i' + \sigma^2 I_{n_i}$ . Finally, when there is interest in predicting the random effects, the SOLUTION (or S) option can be used to request that the estimated BLUPs for the random effects,  $b_i$ , be displayed (in addition to standard errors for predictions based on the expression for  $\text{Var}(\hat{b}_i - b_i)$  given in Section 8.7). Alternatively, the predicted values of the response,  $\hat{Y}_i = X_i \hat{\beta} + Z_i \hat{b}_i$ , can be requested by using the OUTPRED (or OUTP) option on the MODEL statement. This option specifies a SAS data-set than contains the predicted values of  $Y_i$ , denoted by the variable name Pred, and some related quantities. For example, to obtain the estimated BLUPs for the random effects and predicted values of the response,  $Y_i$ , we can use the illustrative SAS commands given in [Table 8.16](#). The

OUTPRED option specifies an output SAS data-set containing the predicted values,  $\hat{Y}_i$ , whereas the SOLUTION option on the RANDOM statement requests that the estimated BLUPs be produced as part of the standard output from PROC MIXED. Inclusion of the Output Delivery System (ODS) statement creates a SAS data-set containing the estimated BLUPs. Predicted values of the outcome, at occasions other than those actually observed, can also be obtained by including “pseudo-observations” in the data set that have missing values for the outcome variable and the desired values of the covariates.

**Table 8.16** Illustrative commands for obtaining the estimated BLUPs and the predicted responses from a linear mixed effects model, with randomly varying intercepts and slopes, using PROC MIXED in SAS.

---

```
ODS OUTPUT SOLUTIONR=bluptable;
PROC MIXED;
  CLASS id group;
  MODEL y=group time group*time/S CHISQ OUTPRED=yhat;
  RANDOM INTERCEPT time/TYPE=UN SUBJECT=id SOLUTION G V;
  PROC PRINT DATA=yhat;
    VAR id group time y Pred;
  PROC PRINT DATA=bluptable;
```

---

The alert reader would have noticed that the residual error variance,  $\sigma^2$ , has not been included on the RANDOM statement. Instead, it is included in an implicit REPEATED statement. Recall that the repeated statement is used to specify assumptions about the nature of the covariance among the errors. When the REPEATED statement is not included in PROC MIXED, it is assumed, by default, that the covariance among the errors,  $R_i = \sigma^2 I_{n_i}$ . To fit hybrid models that include both random effects and correlated errors, it is necessary to include both the RANDOM statement and the REPEATED statement. For example, to fit a hybrid model with (1) randomly varying intercepts and slopes, (2) within-subject errors with an exponential covariance structure, and (3) independent measurement or sampling errors, we can use the illustrative SAS commands given in [Table 8.17](#). On the REPEATED statement we use the option TYPE=SP(EXP)(time) to specify an exponential covariance structure for the within-subject errors that depends on time. This command exploits the spatial covariance structures option built into PROC MIXED. Finally, the option LOCAL requests that a diagonal matrix,  $\sigma^2 I_{n_i}$ , be added to the exponential covariance structure for  $R_i$ .

**Table 8.17** Illustrative commands for a linear mixed effects model, with randomly varying intercepts and slopes, within-subject errors with an exponential covariance, and independent measurement errors using PROC MIXED in SAS.

---

```
PROC MIXED;
  CLASS id group;
  MODEL y=group time group*time/S CHISQ;
  REPEATED/TYPE=SP(EXP)(time) LOCAL SUBJECT=id;
  RANDOM INTERCEPT time/TYPE=UN SUBJECT=id G V;
```

---

## 8.10 FURTHER READING

Useful reviews of the linear mixed effects models, targeted at applied researchers, can be found in the articles by Feldman (1988), Gibbons et al. (1988), Naumova et al. (2001), and Chapters 3 and 4 of Singer and Willett (2003). A comprehensive, but more mathematically challenging discussion of linear mixed effects models can be found in Chapter 3 of Verbeke and Molenberghs (2000) and in the review article by Cnaan *et al.* (1997).

An excellent, non-technical, discussion of the notion of shrinkage can be found in Efron and Morris (1977); also, see the discussion of prediction of random effects in Naumova et al. (2001).

Finally, a tutorial description of fitting linear mixed effects models using PROC MIXED in SAS can be found in Singer (1998); also see Chapters 6 and 7 of Littell et al. (1996) and Chapter 8 of Verbeke and Molenberghs (2000).

# Bibliographic Notes

Harville (1977) introduced a general class of linear mixed effects models suitable for the analysis of repeated measures and growth curves; also see Hartley and Rao (1967). The idea of allowing certain regression coefficients to vary randomly across individuals was also a recurring theme in the early contributions to growth curve analysis by Wishart (1938), Box (1950), Rao (1958), Potthoff and Roy (1964), and Grizzle and Allen (1969); these early contributions to growth curve modeling laid the foundation for the linear mixed effects model. Laird and Ware (1982), Jennrich and Schluchter (1986), Laird et al. (1987), Lindstrom and Bates (1988), Diggle (1988), Chi and Reinsel (1989), and others, drew upon this family to propose a general class of models for longitudinal data. Ware (1985) provides a general overview of the application of linear mixed effects models to repeated measures and longitudinal data; also see Chapter 3 of Davidian and Giltinan (1995) for a concise review of linear mixed effects models for repeated measures data.

The notion of shrinkage was first introduced in a seminal paper by Stein (1955). Best linear unbiased prediction (BLUP) is discussed in Henderson (1963); see Robinson (1991) for an interesting review of the prediction of random effects.

## Problems

**8.1** In a study of exercise therapies, 37 patients were assigned to one of two weightlifting programs (Freund *et al.*, 1986). In the first program (treatment 1), the number of repetitions was increased as subjects became stronger. In the second program (treatment 2), the number of repetitions was fixed but the amount of weight was increased as subjects became stronger. Measures of strength were taken at baseline (day 0), and on days 2, 4, 6, 8, 10, and 12.

The raw data are stored in an external file: `exercise.dat`

Each row of the data set contains the following nine variables:

ID Treatment  $Y_1$   $Y_2$   $Y_3$   $Y_4$   $Y_5$   $Y_6$   $Y_7$

*Note:* The categorical variable Treatment is coded 1 = Program 1 (increase number of repetitions), 2 = Program 2 (increase amount of weight).

**8.1.1** On a single graph, construct a time plot that displays the mean strength versus time (in days) for the two treatment groups. Describe the general characteristics of the time trends for the two exercise programs.

**8.1.2** Read the data from the external file and put the data in a “univariate” or “long” format, with 7 “records” per patient.

**8.1.3** Fit a model with randomly varying intercepts and slopes, and allow the mean values of the intercept and slope to depend on treatment group (i.e., include main effect of treatment, a linear time trend, and a treatment by linear time trend interaction as fixed effects).

(a) What is the estimated variance of the random intercepts?

(b) What is the estimated variance of the random slopes?

(c) What is the estimated correlation between the random intercepts and slopes?

(d) Give an interpretation to the magnitude of the estimated variance of the random intercepts.

For example, “approximately 95% of subjects have baseline measures of strength between a and b” (calculate the limits of the interval between a and b).

(e) Give an interpretation to the magnitude of the estimated variance of the random slopes.

**8.1.4** Is a model with only randomly varying intercepts defensible? Explain?

**8.1.5** What are the mean intercept and slope in the two exercise programs?

**8.1.6** Based on the previous analysis, interpret the effect of treatment on changes in strength. Does your analysis suggest a difference between the two groups?

**8.1.7** What is the estimate of  $\text{Var}(Y_{i1}|b_i)$ ? What is the estimate of  $\text{Var}(Y_{il})$ ? Explain the difference.

**8.1.8** Obtain the predicted (empirical BLUP) intercept and slope for each subject.

**8.1.9** Using any standard linear regression procedure, obtain the ordinary least squares (OLS) estimates of the intercept and slope from the regression of strength on time (in days) for subject 24 (ID = 24). That is, restrict the analysis to data on subject 24 only and estimate that subject's intercept and slope.

**8.1.10** For subject 24 (ID = 24), compare the predicted intercepts and slopes obtained in Problems 8.1.8 and 8.1.9. How and why might these differ?

**8.2** AIDS Clinical Trial Group (ACTG) study 193A was a randomized, double-blind, study of AIDS patients with advanced immune suppression (CD4 counts of  $\leq 50$  cells/mm $^3$ ) (Henry et al., 1998). Patients were randomized to one of four daily regimens containing 600 mg of zidovudine:

- (1) zidovudine alternating monthly with 400 mg didanosine;
- (2) zidovudine plus 2.25 mg of zalcitabine;
- (3) zidovudine plus 400 mg of didanosine;
- (4) zidovudine plus 400 mg of didanosine plus 400 mg of nevirapine.

In the analyses reported in Section 8.8, the first three treatment groups were combined and compared to the fourth.

Measurements of CD4 counts were scheduled to be collected at baseline and at 8-week intervals during follow-up. However, the CD4 count data are unbalanced due to mistimed measurements and missing data that resulted from skipped visits and dropout. The number of measurements of CD4 counts during the first 48 weeks of follow-up varied from 1 to 9, with a median of 4. CD4 count refers to the number of T-lymphocyte cells in the body; these cells are directly affected by the HIV virus. A normal CD4 count is approximately 800 to 1000; a CD4 count below 200 is one of the diagnostic criteria for AIDS established by the Centers for Disease Control and Prevention (CDC).

The raw data are stored in an external file: `cd4.dat`

Each row of the data set contains the following six variables:

ID Group Age Gender Week Log(CD4 + 1)

*Note:* The categorical variable Group is coded 1 = zidovudine alternating monthly with 400 mg didanosine, 2 = zidovudine plus 2.25 mg of zalcitabine, 3 = zidovudine plus 400 mg of didanosine, and 4 = zidovudine plus 400 mg of didanosine plus 400 mg of nevirapine. The variable Week represents time since baseline (in weeks).

**8.2.1** On a single graph, construct a smoothed time plot that displays the mean log CD4 counts versus time (in weeks) for the four treatment groups. Describe the general characteristics of the time trends for the four groups.

**8.2.2** Fit a model where each patient's response trajectory is represented by a randomly varying piecewise linear spline with a knot at week 16. That is, fit a model with random intercepts and two randomly varying slopes, one slope for the changes in log CD4 counts before week 16, another slope for the changes in response after week 16. Allow the average slopes for changes in response before and after week 16 to vary by group, but assume the mean response at baseline is the same in the four groups.

**8.2.3** Is a model with only randomly varying intercepts defensible? Explain?

**8.2.4** Construct a 6-degrees-of-freedom test of the null hypothesis of no treatment group differences in the changes in log CD4 counts.

**8.2.5** Based on the previous analysis, interpret the effects of treatment on changes in log CD4 counts from baseline to week 16.

**8.2.6** Based on the previous analysis, interpret the effects of treatment on changes in log CD4 counts from week 16 to week 40.

**8.2.7** Using the estimates of the fixed effects from the previous analysis, construct a time plot that displays the *estimated* mean log CD4 counts versus time (in weeks) for the four treatment groups. Does the plot suggest that one treatment regimen is superior to the others in terms of short-term (40 weeks) changes in CD4 counts?

<sup>1</sup> Note that, in this section, in a slight abuse of notation, we are using  $\beta$  to denote a fixed parameter and  $\beta_i$  to denote a random variable.

<sup>2</sup> Alternatively, this constraint can be relaxed by assuming that some components of  $\beta_i$  are constant, not random. However, when some components of  $\beta_i$  are assumed to be constant, then the repeated measures on each individual no longer follow a regression model with distinct regression coefficients for each individual in stage 1.

<sup>3</sup> It is difficult to attribute the popularization of the so-called NIH method to any single biostatistician at NIH. During their time at NIH, Sam Greenhouse, Max Halperin, and Jerry Cornfield introduced many biostatisticians to this technique.

<sup>4</sup> Z statistics are intentionally omitted from [Table 8.3](#) because, in general, we do not recommend testing hypotheses about the covariance parameters using Wald tests. In particular, the sampling distribution of a variance parameter estimate does not have an approximate normal distribution when the sample size is relatively small and the population variance is close to zero (see Chapter 4, Section 4.4).

<sup>†</sup> This section provides a more technical presentation of prediction of random effects and the notion of shrinkage; it can be omitted at first reading without loss of continuity.

# *Chapter 9*

## *Fixed Effects versus Random Effects Models*

### **9.1 INTRODUCTION**

In the previous chapter we discussed linear mixed effects models for longitudinal data. The key feature of the models was the introduction of random effects to account for natural heterogeneity in the study population. An alternative, but closely related, class of linear regression models for longitudinal data have their origins in the econometrics literature where they are known as “fixed effects” models. These so-called fixed effects models are increasingly used in the social sciences, broadly defined, and differ in a number of important ways from the mixed effects models considered in Chapter 8. In this chapter, we review the main features of linear fixed effects models for longitudinal data and discuss their potential advantages and disadvantages relative to linear mixed effects models. We highlight the key differences between the two modeling approaches using a numerical illustration. We also apply both modeling approaches to lung function growth data from the Six Cities Study of Air Pollution and Health. Finally, we discuss a mixed effects model for longitudinal data that incorporates the most desirable features of both classes of models via an appropriate decomposition of between- and within-subject effects.

## 9.2 LINEAR FIXED EFFECTS MODELS

Recall that one important motivation for the use of regression models is the control of confounding. In longitudinal studies where randomization cannot be employed, great emphasis must be placed on the measurement and control of important confounding variables. By construction, regression models allow the assessment of the effects of the covariates of main scientific interest while statistically adjusting or controlling for confounding variables that have also been included in the analysis. Of course, this type of adjustment for confounding has at least two limitations. First, no matter how many confounding variables have been incorporated into the regression model, the resulting analysis is always open to the critique that some critical confounders have been omitted. Second, even if it were possible to have complete consensus on the full list of potential confounders, it is almost certain that some of these variables would be inherently difficult or prohibitively expensive to measure. As a result it is almost inevitable that most non-randomized studies will fail to measure all of the key potential confounders and the results of any regression analysis must be interpreted with some caution. Fixed effects models were developed with the goal of overcoming both of these limitations, at least for one particular type of confounding variable in a longitudinal study.

Linear fixed effects models were introduced by econometricians with the intended goal of eliminating an important potential source of bias from regression models for longitudinal data. The fundamental idea underlying fixed effects models is the control of all potential confounding variables that remain stable across repeated measurement occasions and whose effects on the response are assumed to be constant over time. That is, fixed effects models were developed with the goal of removing the potential confounding effects of both observed and unobserved *time-invariant* confounders from longitudinal analyses, under the assumption that the effects of these confounders on the response remain constant over time. For their application to longitudinal data, fixed effects models require two features of the data: (1) two or more repeated measures of the response variable, and (2) values of the covariates of main interest must vary over measurement occasions, for at least some subset of the sample. The first requirement is trivial and is met, by definition, in all longitudinal studies. The second requirement implies that fixed effects models will be most useful in those settings where the main covariates of scientific interest are time-varying. Conversely, fixed effects models are not useful when it is also of interest to estimate the effects of time-invariant covariates.

Next we consider the statistical formulation of the linear fixed effects model. The notation we use is very similar to that employed in Chapter 8. To accommodate unbalanced data, we assume that there are  $n_i$  repeated measurements of the response on the  $i^{th}$  subject and that each  $Y_{ij}$  is observed at time  $t_{ij}$ . Associated with each response,  $Y_{ij}$ , there is a  $p \times 1$  vector of covariates. The vector of covariates can be partitioned into two main types of covariates: covariates whose values do not change throughout the duration of the study and covariates whose values change over time. The former are referred to as time-invariant or between-subject covariates (e.g., gender and fixed experimental treatments), while the latter are referred to as time-varying or within-subject covariates (e.g., time since baseline, current smoking status, and environmental exposures). In a slight departure from the notation used in previous chapters, we let  $X_{ij}$  denote the  $q \times 1$  vector of time-varying covariates and  $W_{ij}$  denote the  $(p - q) \times 1$  vector of time-invariant covariates. For the latter, the same values of the covariates are replicated in the corresponding rows of  $W_{ij}$ , for  $j = 1, \dots, n_i$ ; so we can drop the second subscript and denote the time-invariant covariates by  $W_i$ . The linear fixed effects model is given by

$$(9.1) \quad Y_{ij} = X'_{ij}\beta + W'_i\gamma + \alpha_i + \epsilon_{ij},$$

where the  $\alpha_i$  are *fixed effects* representing stable (i.e., time-invariant) characteristics of individuals that are not otherwise accounted for by the inclusion of the time-invariant covariates,  $W_i$ , in the model. The model is completed by assuming the random within-subject errors,  $\epsilon_{ij} \sim N(0, \sigma^2_\epsilon)$ . At first glance, this model bears remarkable resemblance to the linear mixed effects model considered in

Chapter 8, specifically, the random intercepts model. One of the key distinctions, however, is that in the fixed effects model formulation the  $\alpha_i$  are considered to be fixed effects, whereas in the linear mixed effects model formulation the  $\alpha_i$  are considered to be random. Unfortunately, there are many conflicting definitions of the terms “fixed effects” and “random effects” in the statistical literature and the use of this nomenclature leads to much confusion among statisticians and practitioners alike. As we will see, the key distinction between *fixed effects* models and *random effects* models is whether the  $\alpha_i$  in (9.1) are assumed to be correlated with the covariates. Later we discuss the similarities and differences between these two classes of models in much greater detail.

We want to compare and contrast the fixed and mixed effects models in terms of their model assumptions. The fixed effects model makes the following assumptions about the relationships between  $X_{ij}$ ,  $W_i$ ,  $\alpha_i$ , and  $\epsilon_{ij}$ . First,  $X_{ij}$  is assumed to be completely independent of the random errors, not only  $\epsilon_{ij}$  but also  $\epsilon_{ij}'$  for  $j' \neq j$ . Importantly, this assumption implies that the current value of the response variable,  $Y_{ij}$ , given  $X_{ij}$ , does not predict the subsequent value of  $X_{i,j+1}$ . **Econometricians refer to this as the assumption that  $X_{ij}$  is strictly exogenous.** Some of the implications of the assumption of strict exogeneity are discussed in Chapter 13, Section 13.5. Second, the fixed effects model allows the  $\alpha_i$  to be correlated with  $X_{ij}$ . It is the latter assumption that sets the fixed effects model apart from the mixed effects model considered in Chapter 8. That is, in the linear mixed effects model we assume that  $X_{ij}$  is *strictly exogenous* but make the additional assumption that the  $\alpha_i$ , now considered random rather than fixed effects (and denoted by  $b_i$  in Chapter 8), are independent of  $X_{ij}$  (and independent of  $W_i$ , and  $\epsilon_{ij}$ ). Thus the mixed effects model, unlike the fixed effects model, requires the additional assumption that the random subject effects are uncorrelated with  $X_{ij}$  at all occasions  $j = 1, \dots, n_i$ . Later we will discuss the implications of these assumptions for inferences about  $\beta$  (and  $\gamma$ ).

At this point it is instructive to examine some features of the fixed effects model in the simplest possible longitudinal design with only two repeated measures. Also for simplicity, we assume that  $X_{ij}$  and  $W_i$  are both scalar, each composed of a single time-varying and time-invariant covariate respectively. The fixed effects model is given by

$$Y_{ij} = \beta X_{ij} + \gamma W_i + \alpha_i + \epsilon_{ij}, \text{ for } j = 1, 2$$

(for simplicity, a model with intercept equal to zero is assumed). Under the assumptions of the fixed effects model,  $\epsilon_{ij} \sim N(0, \sigma^2_\epsilon)$ . It is tempting to try to estimate  $\beta$ ,  $\gamma$ , and the  $\alpha_i$  via ordinary least squares (OLS) regression. However, any attempt to jointly estimate all of these effects will run afoul of the following difficulty: the  $\alpha_i$  are perfectly collinear with  $\gamma W_i$ . Consequently we cannot estimate both  $\alpha_i$  and  $\gamma$  from the data at hand. Thus a very notable feature of fixed effects models is that they provide estimates only of the regression parameters for time-varying covariates. The effects of time-invariant covariates simply cannot be estimated in the fixed effects model formulation because of their perfect collinearity with  $\alpha_i$ .

To estimate  $\beta$ , we can consider the model for the within-subject changes in the response,  $Y_{i2} - Y_{i1}$ ,

$$\begin{aligned} Y_{i2} - Y_{i1} &= \beta X_{i2} + \gamma W_i + \alpha_i + \epsilon_{i2} - (\beta X_{i1} + \gamma W_i + \alpha_i + \epsilon_{i1}) \\ &= \beta(X_{i2} - X_{i1}) + (\epsilon_{i2} - \epsilon_{i1}). \end{aligned}$$

The regression model for the within-subject changes meets all of the usual assumptions of standard linear regression. Thus  $\beta$  can be estimated from the OLS regression of  $(Y_{i2} - Y_{i1})$  on  $(X_{i2} - X_{i1})$ , using any standard linear regression software, because the error terms,  $\epsilon_{i2} - \epsilon_{i1}$ , have mean zero,  $E(\epsilon_{i2} - \epsilon_{i1}) = 0$ , and constant variance,  $\text{Var}(\epsilon_{i2} - \epsilon_{i1}) = 2\sigma^2_\epsilon$ . Note that in the construction of a model for the within-subject changes, both the time-invariant covariate effect and the stable (or time-invariant) characteristics of an individual, denoted by  $\alpha_i$ , have disappeared from the model. This makes it clear that the  $\alpha_i$  cannot possibly have any effect on the estimation of  $\beta$ , regardless of whether or not they are correlated with  $X_{ij}$ . Therefore, if the  $\alpha_i$  are thought of as all those unmeasured, but time-invariant, characteristics of an individual, then the so-called fixed effects estimate of the effect of  $X_{ij}$  on  $Y_{ij}$  is

unbiased even if some of those characteristics are considered to be confounders of the relationship between  $X_{ij}$  and  $Y_{ij}$ . It is in that sense that it can be said that the fixed effects model removes the potential for bias due to confounding by all measured and unmeasured time-invariant characteristics of individuals (the latter are encapsulated in the  $\alpha_i$ ). However, there is one important additional assumption that must not be overlooked. The fixed effects model can only remove the potential confounding by those measured and unmeasured time-invariant covariates *whose effects on the response remain constant over time*. That is, conditional on  $X_{ij}$  and  $W_i$ , it must be assumed that the effect of any time-invariant confounder on  $Y_{i1}$  is the same as on  $Y_{i2}$ .

Although in the illustration above we considered only two repeated measures, the same logic and rationale concerning properties of the fixed effects model apply more generally to the case of more than two repeated measures. In the more general case it can be shown that the fixed effects model estimator of  $\beta$  is the standard OLS estimator for “mean-centered” transformations<sup>1</sup> of  $Y_{ij}$  and the covariates. Here mean-centering simply involves taking averages over repeated measurement occasions separately for each individual and then transforming the response and covariates by subtraction of the subject-specific means for the response and covariates. Specifically, let  $Y_{ij}^* = Y_{ij} - \bar{Y}_i$ , where  $\bar{Y}_i = \frac{1}{n_i} \sum_{j=1}^{n_i} Y_{ij}$ , denote the mean-centered response for the  $i^{th}$  individual at the  $j^{th}$  occasion. In a similar fashion we can define the mean-centered covariates,  $X_{ijk}^* = X_{ijk} - \bar{X}_{ik}$ , where  $\bar{X}_{ik} = \frac{1}{n_i} \sum_{j=1}^{n_i} X_{ijk}$ . Note that mean-centering of  $W_i$ , the time-invariant covariates, removes all of the variation in these covariates so that  $W_i^* = 0$  for all individuals. Because estimation in the fixed effects model is based entirely on within-subject variation in the response and covariates, it is now apparent why the effects of time-invariant covariates cannot be estimated; mean-centering of  $W_i$  sets  $W_i^* = 0$  for all subjects and thereby removes all of the variation in these covariates. In a similar way, mean-centering makes the  $\alpha_i$  disappear from the model.<sup>2</sup> It can be shown that the fixed effects estimator of  $\beta$  is simply the OLS estimator of the regression of the mean-centered response on the mean-centered covariates,

$$Y_{ij}^* = X_{ij}^* \beta + \epsilon_{ij}^*,$$

where  $\epsilon_{ij}^* = \epsilon_{ij} - \bar{\epsilon}_i$ . Interestingly, although the errors  $\epsilon_{ij}^*$  in the mean-centered model are no longer uncorrelated over time for a given individual, the correlation among these errors can be ignored for estimation of  $\beta$ ; that is, for the mean-centered model, the generalized least squares (GLS) estimator of  $\beta$  that accounts for the correlation among the  $\epsilon_{ij}^*$  happens to be identical to the OLS estimator of  $\beta$ . The resulting fixed effects estimator of  $\beta$ , denoted  $\hat{\beta}^{(FE)}$ , is

$$\hat{\beta}^{(FE)} = \left( \sum_{i=1}^N \sum_{j=1}^{n_i} X_{ij}^* X_{ij}^{*\prime} \right)^{-1} \left( \sum_{i=1}^N \sum_{j=1}^{n_i} X_{ij}^* Y_{ij}^* \right).$$

An important property of this estimator is that it is unbiased for  $\beta$ , even in cases where the  $\alpha_i$  are correlated with  $X_{ij}$ .

To summarize, the main appeal of fixed effects models is their statistical control of time-invariant characteristics of individuals whose effects on the response are assumed to remain constant over time. An important feature of these models is that they have the potential to remove bias when there are unmeasured, but stable, characteristics of the subjects that are correlated with time-varying covariates of main scientific interest. Of note, fixed effects models cannot estimate the effects of time-invariant covariates; all such covariates are effectively removed from the analysis, as is their influence on the associations between the time-varying covariates and the response. However, fixed effects models can estimate interactions between time-invariant and time-varying covariates. For example, although the effect of a time-invariant treatment on the response at any particular occasion cannot be estimated, by including an interaction between treatment and time since baseline we can estimate how treatment effect at a particular occasion differs from that at baseline.

## 9.3 FIXED EFFECTS VERSUS RANDOM EFFECTS: BIAS-VARIANCE TRADE-OFF

As we noted earlier, fixed effects models remove the potential for bias due to certain types of confounding variables, namely measured and unmeasured time-invariant confounders whose effects on the response can be assumed to be constant across measurement occasions. This property is not shared by linear mixed effects models. The linear mixed effects model implicitly makes stronger assumptions about the stable, time-invariant, characteristics of individuals. Specifically, mixed effects models treat the  $\alpha_i$  as random, rather than fixed, and assume that they are uncorrelated with the measured covariates included in the regression model; that is, the  $\alpha_i$  are assumed to be uncorrelated with  $X_{ij}$  (and also uncorrelated with  $W_i$  and with  $\varepsilon_{ij}$ ). When these assumptions do not hold for the data at hand, and some of the between-subject variation in the  $\alpha_i$  includes unmeasured characteristics of subjects that are correlated with the covariates of interest,  $X_{ij}$ , then the linear mixed effects model can yield biased estimates of the effects of  $X_{ij}$  on  $Y_{ij}$ . In contrast, the fixed effects model makes milder assumptions about the stable characteristics of individuals, implicitly allowing for the possibility of correlation between the  $\alpha_i$  and  $X_{ij}$ . In the fixed effects model the effects of  $X_{ij}$  on  $Y_{ij}$  are estimated by effectively ignoring all of the between-subject variation and focusing exclusively on the within-subject variation. By focusing only on within-subject variation, we ensure that the  $\alpha_i$  cannot possibly have any effect on the estimation of  $\beta$ , regardless of whether they are correlated with  $X_{ij}$ . Therefore in the fixed effects model the  $\alpha_i$ , representing measured and unmeasured stable characteristics of individuals, cannot confound the estimation of the effects of  $X_{ij}$  on  $Y_{ij}$ .

Because the fixed effects model has the potential to avoid bias in the estimation of time-varying covariate effects, it may seem that it should always be the method of first choice for longitudinal analysis. However, there are two main reasons why the linear mixed effects model might often be preferred. First, in many longitudinal designs there is scientific interest in the effects of both time-varying and time-invariant covariates. Indeed, in many longitudinal studies the primary covariates of scientific interest are time-invariant, such as fixed treatment or exposure groups, or various background characteristics of individuals (e.g., gender, socioeconomic status). Linear mixed effects models allow for the estimation of the effects of both time-varying and time-invariant covariates; in contrast, fixed effects models cannot estimate the effects of time-invariant covariates. Second, the potential advantage of fixed effects models over mixed effects models in terms of bias in the estimation of the effects of time-varying covariates comes at a price: efficiency. As is quite common in statistical estimation, there is a trade-off between bias and efficiency. Although, under certain conditions, the fixed effects model may provide unbiased estimates of the effects of time-varying covariates, it will in general yield larger standard errors for those estimated effects than those produced by the mixed effects model. The reason is as follows. Recall that a time-varying covariate has two main sources of variation: within-subject variation (i.e., the degree to which values of the covariate vary over time within the same individual) and between-subject variation (i.e., the degree to which values of the covariate vary from one individual to another). In many longitudinal studies the between-subject variation in a time-varying covariate is orders of magnitude greater than the within-subject variation. For example, in a longitudinal study of the effect of body mass index (BMI) on pulmonary function, we might expect to see far greater variation in BMI between individuals than fluctuations in BMI within individuals over the duration of the study. Similarly, in many studies of growth and aging, there may be far greater variation in age between individuals than within individuals over the course of the study. Fixed effects models base estimation exclusively on the within-subject variation and completely ignore the between-subject variation. In contrast, mixed effects models capitalize on both sources of variation, between- and within-subject variation, yielding standard errors that can be substantially smaller. In general, the greater the proportion of variation in a time-varying covariate that is between-subject variation, the larger the difference in the

magnitudes of the standard errors yielded by the fixed and mixed effects models.

Thus the choice between fixed and mixed effects models is very often made on a combination of scientific and statistical grounds. On scientific grounds, the choice between these two classes of models is clear-cut when there is scientific interest in the effects of time-invariant covariates. Fixed effects models simply cannot estimate these effects while mixed effects models can. The choice between these two classes of models for the estimation of the effects of time-varying covariates can be made on statistical grounds but requires the balancing of bias versus precision. In settings where there is substantial concern about confounding by stable characteristics of individuals and/or where there is substantial within-subject variation in the time-varying covariate of interest, the fixed effects model will be preferred. Conversely, in settings where most of the major confounding variables have been measured and included in the analysis and/or where the between-subject variation is orders of magnitude larger than the within-subject variation in the time-varying covariate of interest, the mixed effects model will be preferred.

Finally, it is worth mentioning settings where the two approaches are anticipated to yield very similar estimates of effects.<sup>†</sup> Consider the following simple random effects model, with single covariate  $X_{ij}$ :

$$(9.2) \quad Y_{ij} = \beta_1 + \beta_2 X_{ij} + b_i + \epsilon_{ij},$$

where  $b_i \sim N(0, \sigma_b^2)$ . Although the technical details are omitted here, the ML estimate of  $\beta_2$  in this simple random effects model, denoted  $\hat{\beta}_2^{(RE)}$ , can be shown to be a weighted average,

$$(9.3) \quad \hat{\beta}_2^{(RE)} = (1 - w) \hat{\beta}_2^{(FE)} + w \hat{\beta}_2^{(B)},$$

where  $\hat{\beta}_2^{(B)}$  is obtained from the regression of  $\bar{Y}_i$  on  $\bar{X}_i$  and is based only on the between-subject variation in the response and covariate. Recall that in contrast to  $\hat{\beta}_2^{(B)}$ ,  $\hat{\beta}_2^{(FE)}$  is based only on the within-subject variation in the response and covariate. Thus this simple expression demonstrates that the ML estimate of  $\beta_2$  is based on an optimal weighted combination of the within-subject and between-subject sources of variation in the response and covariate. The actual weight,  $w$ , is determined by the equation,

$$(9.4) \quad w = \frac{(1 - \rho_y)\rho_x}{(1 - \rho_y) + n\rho_y(1 - \rho_x)},$$

where  $\rho_y = \frac{\sigma_b^2}{\sigma_b^2 + \sigma_\epsilon^2}$  is the proportion of variability in the response that is due to between-subject variation, and  $\rho_x$  is the corresponding proportion of variability in the covariate that is due to between-subject variation. From (9.3) it is apparent that when  $w \approx 0$ , we expect  $\hat{\beta}_2^{(RE)} \approx \hat{\beta}_2^{(FE)}$ ; conversely, when  $w \neq 0$ , we expect  $\hat{\beta}_2^{(RE)}$  to differ from  $\hat{\beta}_2^{(FE)}$ , unless  $\hat{\beta}_2^{(FE)} \approx \hat{\beta}_2^{(B)}$ . Further examination of the expression for  $w$  given by (9.4) reveals when  $w$  is expected to be small. First, for a fixed value of  $\rho_x$ ,  $w$  is expected to be small when  $\rho_y$  is large, such as when  $\rho_y \rightarrow 1$ ,  $w \rightarrow 0$ . Thus  $\hat{\beta}_2^{(RE)} \approx \hat{\beta}_2^{(FE)}$  when the between-subject variation in the *response* is large relative to the within-subject variation. For example, the two approaches are anticipated to yield very similar estimates of effects when the pairwise correlation among repeated measures is large and close to one. Put another way, when the within-subject variation in the *response* is relatively small the repeated measures on an individual are highly reliable. Therefore the most precise estimate of  $\beta_2$  is obtained by relying more heavily on how changes in the covariate are associated with changes in the response *within the same subject*. Second,  $w$  is expected to be small when  $\rho_x$  is small, such as when  $\rho_x \rightarrow 0$ ,  $w \rightarrow 0$ . Thus  $\hat{\beta}_2^{(RE)} \approx \hat{\beta}_2^{(FE)}$  when the within-subject variation in the *covariate* is large relative to the between-subject variation in the covariate. For example, in many designed longitudinal studies there are covariates that vary systematically over time but are fixed by design of the study. In many of these studies  $X_{ij} = X_j$  for all subjects, e.g.,  $X_j$  might denote time since baseline when the measurement occasions are fixed by the study design. This implies that all of the variation in the covariate is within-subject variation and  $\rho_x = 0$ . Another closely related example arises in crossover study designs where  $X_{ij}$  is a treatment group indicator denoting the treatment received by the  $i^{th}$  subject at the  $j^{th}$  occasion. In many classical crossover designs, each subject receives all of the treatments being compared but in a random

sequence or order (see Chapter 21). For these classical crossover study designs,  $\bar{X}_i = \bar{X}$  for all subjects, which implies that all of the variation in the covariate is within-subject variation ( $\rho_x = 0$ ). Finally, examination of (9.4) reveals that when the number of repeated measures,  $n$  (here assumed to be equal for all subjects), is large, then  $w$  will tend to be small. So, in longitudinal studies with a very large number of repeated measurements,  $\hat{\beta}_2^{(RE)}$  and  $\hat{\beta}_2^{(FE)}$  are expected to be similar.

## 9.4 RESOLVING THE DILEMMA OF CHOOSING BETWEEN FIXED AND RANDOM EFFECTS MODELS

Fixed and random effects models yield estimators of the effects of time-varying covariates with different desirable properties. Much of the statistical literature comparing these two alternative approaches presents them as mutually exclusive choices; that is, one must choose between the two types of models, making a choice between the potential bias associated with the random effects model formulation and the increased sampling variability associated with the fixed effects model formulation. However, there is a third option, albeit one that does not appear to have been so widely adopted in practice: develop a model that capitalizes on most of the appealing features of both the random and fixed effects model formulations. Specifically, it is possible to specify a linear mixed effects model that explicitly allows, when it is deemed necessary, separate estimation of the effect of a time-varying covariate from the two distinct sources of variation in that covariate: one estimate based exclusively on within-subject variation in the response and covariate, the other estimate based exclusively on between-subject variation in the response and covariate. Indeed, it is possible to calculate these estimates for any subset of the time-varying covariates. Such a model recognizes that longitudinal data can provide two distinct sources of information about the relationship between a time-varying covariate and the response: (1) between-subject information reflected in the fact that, at any occasion, different individuals have different values of the covariate and response, and (2) within-subject information reflected in the fact that the covariate and response change over time within subjects. These can be thought of as “cross-sectional” and “longitudinal” information, respectively. Such a model also recognizes that these two sources of information can potentially provide conflicting signals about the nature and magnitude of the covariate effect; the latter is highlighted with a graphical illustration in Section 9.5.

In the standard linear mixed effects model, the estimated effect of a time-varying covariate is based on an optimal combination of the within-subject and between-subject variation. However, when the longitudinal and cross-sectional information are in conflict, these two sources of information should not be combined. Conversely, when they are not in conflict, it is advantageous to optimally combine them to yield the most precise estimate of the covariate effect. These are the principles that guide the specification of a linear mixed effects model that explicitly allows, when necessary, for separate estimation of the effects of time-varying covariates based exclusively on the within- or the between-subject variation in the response and covariates. That is, the model makes an appropriate decomposition of between- and within-subject effects. The case for such a model formulation is at least fivefold:

1. The model allows for joint estimation of the effects of both time-varying and time-invariant covariates, thereby overcoming an important limitation of the fixed effects model.
2. By including a vector of random effects, the model allows for heterogeneity in the effects of certain time-varying covariates and a more flexible model for the marginal covariance among the repeated measures of the response.
3. The model allows separate estimation of the effects of time-varying covariates based on within- and between-subject variation. The resulting estimator based on the within-subject variation shares all of the desirable properties of the fixed effects estimator, i.e., the estimator is unbiased even when  $b_i$  is correlated with a time-varying covariate.
4. The model allows for a formal statistical comparison of the estimates based on the within-subject and between-subject variation. Moreover, when there is sufficient evidence that they are not discernibly different, a combined estimate based on both sources of variation can be obtained; the resulting estimator shares the desirable property with the random effects estimator of being more efficient than the fixed effects estimator.

5. Finally, maximum likelihood estimation of the covariate effects within a linear mixed effects model yields valid inferences when data are missing at random (MAR), but not necessarily missing completely at random (MCAR); see Section 4.3 for the definitions of, and the distinction between, MCAR and MAR. In contrast, OLS estimation of fixed effects models can produce biased estimates of covariate effects when data are MAR.

This approach, combining the appealing features of both the fixed and random effects formulation without any of the major limitations of either approach, provides a more attractive option than having to choose between adopting either a fixed or random effects model.

Next we outline the specification of this model. Following the notation introduced earlier, we let  $X_{ij}$  denote the  $q \times 1$  vector of time-varying covariates and  $W_i$  denote the  $(p - q) \times 1$  vector of time-invariant covariates. We denote the mean-centered covariates by  $X_{ijk}^*$ , where  $X_{ijk}^* = X_{ijk} - \bar{X}_{ib}$ , for  $k = 1, \dots, q$ . The mean-centered covariates can be grouped together in a  $q \times 1$  vector denoted by  $\mathbf{x}_{ij}^*$ ; similarly the means of the covariates can be grouped together in a  $q \times 1$  vector denoted by  $\bar{\mathbf{X}}_i$ . The linear mixed effects model that allows for simultaneously estimating cross-sectional and longitudinal effects of  $X_{ij}$  on  $Y_{ij}$  is then given by

$$(9.5) \quad Y_{ij} = \mathbf{x}_{ij}^{*\prime} \beta^{(L)} + \bar{\mathbf{X}}_i' \beta^{(C)} + W_i' \gamma + Z_{ij}' b_i + \epsilon_{ij},$$

where  $b_i$  is a vector of random effects, with  $b_i \sim N(0, G)$ ,  $Z_{ij}$  is the design vector for the random effects and is a subset of the components of  $X_{ij}$ , and  $\epsilon_{ij} \sim N(0, \sigma_\epsilon^2)$ . Note that the inclusion of a vector of random effects,  $b_i$ , allows for heterogeneity in the effects of certain time-varying covariates and a flexible model for the marginal covariance among the repeated measures of the response (see Section 8.3).

This model allows both cross-sectional effects,  $\beta^{(C)}$ , and longitudinal effects,  $\beta^{(L)}$ , to be modeled simultaneously. The interpretation of the model parameters,  $\beta^{(C)}$  and  $\beta^{(L)}$ , becomes more transparent when the implied models for the average response (where averaging is over time) and within-subject changes in the response are considered. First, consider the model for the time-averaged response,

$$\bar{Y}_i = \frac{1}{n_i} \sum_{j=1}^{n_i} Y_{ij}$$

$$\bar{Y}_i = \bar{\mathbf{X}}_i' \beta^{(C)} + W_i' \gamma + e_i,$$

where  $e_i = \frac{1}{n_i} \sum_{j=1}^{n_i} (Z_{ij}' b_i + \epsilon_{ij}) = \bar{Z}_i' b_i + \bar{\epsilon}_i$ . The regression parameters  $\beta^{(C)}$  represent the cross-sectional effects of  $X_{ij}$  when different individuals having distinct average values for  $X_{ij}$  are compared and contrasted. Next consider the model for within-subject changes, say  $Y_{ij} - Y_{i1}$ . The model for the within-subject changes is given by

$$\begin{aligned} Y_{ij} - Y_{i1} &= \mathbf{x}_{ij}^{*\prime} \beta^{(L)} + \bar{\mathbf{X}}_i' \beta^{(C)} + W_i' \gamma + Z_{ij}' b_i + \epsilon_{ij} \\ &\quad - (X_{i1}^{*\prime} \beta^{(L)} + \bar{\mathbf{X}}_i' \beta^{(C)} + W_i' \gamma + Z_{i1}' b_i + \epsilon_{i1}) \\ &= (X_{ij}' - X_{i1}') \beta^{(L)} + (Z_{ij}' - Z_{i1}') b_i + (\epsilon_{ij} - \epsilon_{i1}) \\ &= (X_{ij}' - X_{i1}') \beta^{(L)} + e_{ij}, \end{aligned}$$

where  $e_{ij} = (Z_{ij}' - Z_{i1}') b_i + (\epsilon_{ij} - \epsilon_{i1})$ . Therefore, in the model for the within-subject changes,  $\beta^{(L)}$  represents a vector of regression parameters for longitudinal effects of  $X_{ij}$ , describing how within-subject changes in the covariates are related to within-subject changes in the response. Note that when it is assumed that the cross-sectional and longitudinal effects of  $X_{ij}$  are the same, so that  $\beta^{(C)} = \beta^{(L)} = \beta$ , the model given by (9.5) collapses to the standard linear mixed effects model,

$$\begin{aligned} Y_{ij} &= \mathbf{x}_{ij}^{*\prime} \beta^{(L)} + \bar{\mathbf{X}}_i' \beta^{(C)} + W_i' \gamma + Z_{ij}' b_i + \epsilon_{ij} \\ &= \mathbf{x}_{ij}^{*\prime} \beta + \bar{\mathbf{X}}_i' \beta + W_i' \gamma + Z_{ij}' b_i + \epsilon_{ij} \\ &= (X_{ij}' \beta - \bar{\mathbf{X}}_i' \beta) + \bar{\mathbf{X}}_i' \beta + W_i' \gamma + Z_{ij}' b_i + \epsilon_{ij} \\ &= X_{ij}' \beta + W_i' \gamma + Z_{ij}' b_i + \epsilon_{ij}. \end{aligned}$$

One important advantage of simultaneously modeling cross-sectional and longitudinal effects is that formal comparisons can be made by testing  $H_0: \beta^{(C)} = \beta^{(L)}$  (or by comparing certain components of  $\beta^{(C)}$  with the corresponding components of  $\beta^{(L)}$ ). In the econometrics literature, the comparison of the estimates of the effects of time-varying covariates from fixed and random effects models is

known as the “Hausman test” (Hausman, 1976). The Hausman test is often presented as a classical test of whether the fixed or random effects model should be used. However, as noted earlier, in many longitudinal designs there is scientific interest in both the effects of time-varying covariates and time-invariant covariates (e.g., fixed treatment or exposure groups). In these settings the fixed effects model is not appealing because it precludes the estimation of the effects of time-invariant covariates. On the other hand, it is of interest to compare the cross-sectional and longitudinal effects of time-varying covariates because the longitudinal effects are not prone to confounding by measured and unmeasured stable characteristics of individuals.

The testing of  $H_0: \beta^{(C)} = \beta^{(L)}$  in the linear mixed effects model is very similar in spirit to the Hausman test. However, within this modeling framework we can also consider a much broader range of tests, such as tests of equality of subsets of the components of  $\beta^{(C)}$  and  $\beta^{(L)}$ . In cases where there is insufficient evidence that certain components of  $\beta^{(C)}$  and  $\beta^{(L)}$  differ, we can obtain a single estimate from the linear mixed effects model that is based on the optimal combination of the between-subject (or cross-sectional) and within-subject (or longitudinal) sources of variation. Moreover the model allows for estimation of the effects of time-invariant covariates,  $\gamma$ , regardless of whether it is assumed that  $\beta^{(C)} = \beta^{(L)}$

## 9.5 LONGITUDINAL AND CROSS-SECTIONAL INFORMATION

In this section we highlight how the two sources of information for a time-varying covariate, longitudinal (or within-subject) and cross-sectional (or between-subject), can potentially provide conflicting signals about the nature and magnitude of the covariate effect. To fix ideas, consider a study of aging. In Chapter 1, when we discussed the main distinctions between a longitudinal and cross-sectional study, we emphasized that the assessment of within-subject changes in the response due to aging can only be achieved within a longitudinal study design. In a cross-sectional study, where the response is measured at a single occasion, we cannot estimate the effect of aging (an inherently within-subject effect); instead, we can only make comparisons among sub-populations that happen to differ in age. However, when the effect of aging is determined from a cross-sectional study, there is the danger that it is potentially confounded with cohort effects.

When an initial cross-sectional sample is measured repeatedly through time, it is then possible to make comparisons of longitudinal and cross-sectional estimates of changes in the response. For example, the Muscatine Coronary Risk Factor (MCRF) study enrolled five cohorts of children, initially aged 5–7, 7–9, 9–11, 11–13, and 13–15 years. Repeated measurements of obesity, based on BMI, were obtained biennially, from 1977 to 1981, with the objective of determining whether the risk of obesity increased with age. Note that the data from the MCRF study are unbalanced over time when the age of the child is used as the metamer for time. That is, baseline measurements are taken at the same calendar time (1977) for all subjects but age at entry to the study varies with subjects. As a result there are two potential sources of information about changes in BMI with age. First, there is cross-sectional or between-subject information about how BMI changes with age in the baseline observations obtained in 1977, since children enter the study at different ages. Similar cross-sectional information is also available in 1979 and 1981 (or, equivalently, in the average of the repeated measures over time). Second, longitudinal or within-subject information arises because children are measured repeatedly over time, yielding measurements of BMI at different ages. These two sources of information may provide conflicting information about how BMI changes with age.

As we discussed in previous sections, when a study provides both longitudinal and cross-sectional information, somewhat greater care must be exercised in specifying models for the response to avoid confounding of longitudinal effects with cross-sectional effects. The linear mixed effects model presented in Section 9.4 includes separate parameters that represent the cross-sectional and longitudinal effects of age on the response and allows simultaneous estimation of both types of effects. This makes it possible to compare the cross-sectional and longitudinal effects, and to report separate effects where necessary, or estimate a combined effect, based on both sources of information, if appropriate.

In a small departure from the notation used in the previous section, we let  $\text{Age}_{ij}$  denote the age of the  $i^{th}$  subject at the  $j^{th}$  measurement occasion. A model for aging, with decomposition of between- and within-subject effects, is given by

$$Y_{ij} = \beta_1 + \beta_2^{(L)}(\text{Age}_{ij} - \bar{\text{Age}}_i) + \beta_2^{(C)}\bar{\text{Age}}_i + W_i'\gamma + b_{1i} + b_{2i}\text{Age}_{ij} + \epsilon_{ij},$$

where  $\beta_2^{(C)}$  represents the cross-sectional effect of age since it describes how the mean response at any occasion varies with age at that occasion. In contrast,  $\beta_2^{(L)}$  represents the longitudinal effect of age since it describes how within-subject changes in the response are related to within-subject changes in age. Differences between  $\beta_2^{(C)}$  and  $\beta_2^{(L)}$  can arise when there are cohort or period effects. Cohort effects will introduce bias in the cross-sectional estimate but not the longitudinal estimate. Period effects will introduce bias in the longitudinal estimate but not the cross-sectional estimate. Alternatively, differences between  $\beta_2^{(C)}$  and  $\beta_2^{(L)}$  can be due to the biasing effects of selective dropouts. When  $\beta_2^{(C)} = \beta_2^{(L)} = \beta_2$ , the model given above simplifies to

$$Y_{ij} = \beta_1 + \beta_2\text{Age}_{ij} + W_i'\gamma + b_{1i} + b_{2i}\text{Age}_{ij} + \epsilon_{ij},$$

the standard linear mixed effects model. On the other hand, when  $\beta_2^{(C)} \neq \beta_2^{(L)}$  but the model for the

data does not allow for separate estimation of the cross-sectional and longitudinal effects on the response, then  $\beta_2$  cannot be interpreted as a pure longitudinal effect of aging. Instead,  $\beta_2$  is some weighted combination of  $\beta^{(C)}_2$  and  $\beta^{(L)}_2$ , with weights determined by the relative magnitudes of the between- and within-subject variation (see [equations \(9.3\)](#) and [\(9.4\)](#) in Section 9.3). Such a weighted combination of  $\beta^{(C)}_2$  and  $\beta^{(L)}_2$  may not reflect an effect of subject-matter interest. That is, failure to distinguish between cross-sectional and longitudinal effects can result in a distorted estimate of the effect of age that reflects neither the cross-sectional nor the longitudinal effect of age on the response.

# Illustration

The main distinction between cross-sectional and longitudinal effects is highlighted in the following simple illustration. Suppose that three age-cohorts of children, initially aged 5, 6, and 7 years, are measured at baseline and followed annually for three years. Suppose that the cross-sectional effect of age on the response is linear, with

$$E(\bar{Y}_i) = \beta^{(C)} \bar{\text{Age}}_i$$

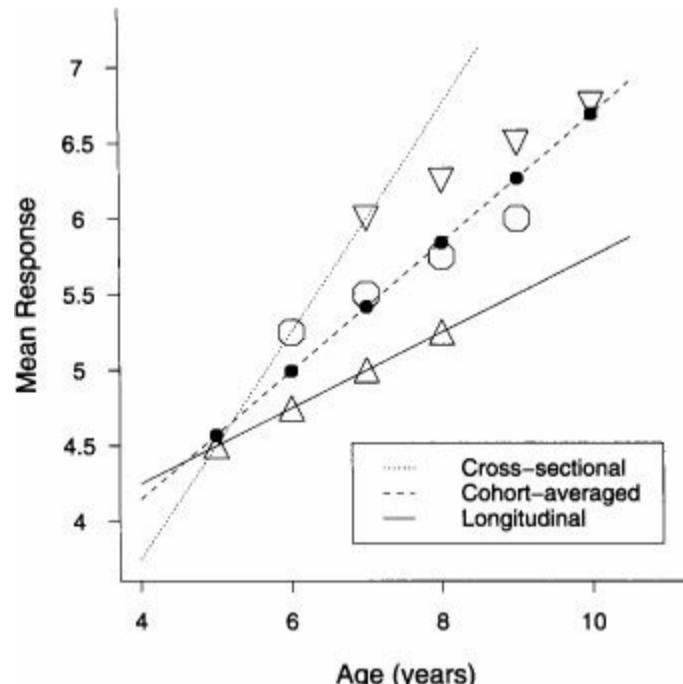
(for simplicity, a model with intercept equal to zero is assumed). The same cross-sectional relationship, with slope  $\beta^{(C)}$ , is assumed to hold at every measurement occasion. The mean response is also assumed to increase linearly with changes in age in each cohort

$$E(Y_{ij} - \bar{Y}_i) = \beta^{(L)} (\text{Age}_{ij} - \bar{\text{Age}}_i),$$

but with slope  $\beta^{(L)} \neq \beta^{(C)}$ . Thus in this model  $\beta^{(C)}$  represents the cross-sectional effect of age, whereas  $\beta^{(L)}$  represents the longitudinal effect of age.

A graphical representation of this model for the mean response versus age, with  $\beta^{(C)} = 0.75$  and  $\beta^{(L)} = 0.25$ , is given in [Figure 9.1](#). In this illustration there is a discernible difference between the longitudinal (solid line) and cross-sectional (dotted line) effects of aging on the mean response. When these differences between the longitudinal and cross-sectional effects of aging are ignored, the observed changes in the mean response with age of measurement (see dashed line in [Figure 9.1](#)) reflect a weighted combination of  $\beta^{(C)}$  and  $\beta^{(L)}$ . Recall from [equation \(9.3\)](#) that the ML estimate from the naive analysis assuming  $\beta^{(C)} - \beta^{(L)} = \beta$  is the following weighted average:

**Fig. 9.1** Plot of the longitudinal, cross-sectional, and cohort-averaged regression lines for the three age-cohorts:  $\Delta$  denotes mean response of the age-cohort of children initially aged 5 years;  $\circ$  denotes mean response of the age-cohort of children initially aged 6 years; and  $\nabla$  denotes mean response of the age-cohort of children initially aged 7 years, ( $\bullet$  denotes mean response when averaged over the three age cohorts).



$$\hat{\beta} = (1 - w) \hat{\beta}^{(L)} + w \hat{\beta}^{(C)},$$

where

$$w = \frac{(1 - \rho_y) \rho_x}{(1 - \rho_y) + n \rho_y (1 - \rho_x)}.$$

Here  $n = 4$  and  $\rho_x$  denotes the proportion of variability in age due to between-subject variation relative to within-subject variation. In this illustration, approximately 35% of the variation is due to between-subject variation and 65% is due to within-subject variation ( $\rho_x = 0.3478$ ). For this study design, with  $n$  and  $\rho_x$  fixed by design, the resulting estimate of  $\beta$  depends on the magnitude of  $\rho_y$ , the correlation among the repeated measures of the response. Specifically, the estimate of  $\beta$  can range from 0.25 (when  $\rho_y = 1$  and  $w = 0$ ) to 0.424 (when  $\rho_y = 0$  and  $w = 0.3478$ ). The dashed line in [Figure 9.1](#) depicts the case where  $\rho_y = 0$ . In general, when differences between the longitudinal and cross-

sectional effects are ignored, the naive analysis assuming  $\beta^{(C)} = \beta^{(L)} = \beta$  estimates a weighted average of  $\beta^{(C)}$  and  $\beta^{(L)}$ .

This simple illustration highlights why standard regression models for longitudinal data that do not incorporate separate cross-sectional and longitudinal effects of time-varying covariates can potentially yield misleading inferences. That is, failure to acknowledge that the cross-sectional effect differs from the longitudinal effect can lead to a conclusion about the effect of the covariate that confounds one effect with the other.

## 9.6 CASE STUDY

Next we illustrate the main ideas presented in this chapter using lung function growth data on a randomly selected subset of female participants from the Six Cities Study of Air Pollution and Health. The random sample consists of 300 girls from Topeka, Kansas, one of the participating cities in the study. Each girl had a minimum of 1 and a maximum of 12 measurements of FEV<sub>1</sub>, height and age over the course of the study. One outlying observation was removed and all analyses are based on the data from 299 girls (with a total of 1993 measurements).

The goal of the analysis is to assess the relationship between FEV<sub>1</sub> and age. There are two sources of information about the relationship between FEV<sub>1</sub> and age. First, there is “cross-sectional” or between-subject information that arises from the comparison of children of different ages. Second, there is “longitudinal” or within-subject information that arises because children are measured repeatedly over time, yielding measurements of FEV<sub>1</sub> at different ages.

We begin with a so-called fixed effects analysis of these data. Recall that the fixed effects analysis focuses exclusively on the “longitudinal” information about the relationship between FEV<sub>1</sub> and age. To adjust for variation in a child’s stature, values for FEV<sub>1</sub> were divided by height squared and log-transformed; this has been shown to be an effective yet parsimonious adjustment for height. We consider the following fixed effects model for log(FEV<sub>1</sub>/height<sup>2</sup>):

$$E(Y_{ij}) = \beta_1 + \beta_2 \text{Age}_{ij} + \alpha_i,$$

where  $Y_{ij}$  is the log(FEV<sub>1</sub>/height<sup>2</sup>) for the  $i^{th}$  child at the  $j^{th}$  occasion,  $\text{Age}_{ij}$  is the age for the  $i^{th}$  child at the  $j^{th}$  occasion, and  $\alpha_i$  are fixed effects representing stable (i.e., time-invariant) characteristics of the children. Because the model includes an intercept,  $\beta_1$ , 298 dummy or indicator variables must be included in the model to estimate the  $\alpha_i$ . In this model,  $\beta_2$  is the “longitudinal” effect of age. Estimates of  $\beta_1$ ,  $\beta_2$ , and  $\alpha_i$  can be obtained, via ordinary least squares (OLS), using any standard linear regression software. The results from fitting this fixed effects model are reported in [Table 9.1](#). These results indicate that there is a significant effect of age. The estimated coefficient for age, 0.0298, indicates that each one year increase in age is associated with an  $e^{0.0298} = 1.030$  or 3.0% relative increase in FEV<sub>1</sub> (adjusted for height<sup>2</sup>). Also reported in [Table 9.1](#) are the estimates of the subject-specific effects ( $\beta_1 + \alpha_i$ ) for the first three subjects in the sample.

**Table 9.1** Regression estimates and standard errors for the fixed effects model for the log(FEV<sub>1</sub>/height<sup>2</sup>) data from the Six Cities Study.

Variable	Estimate	SE	Z
ID 1	-0.399	0.0252	-15.81
ID 2	-0.302	0.0238	-12.72
ID 3	-0.243	0.0223	-10.91
:	:	:	:
Age ( $\times 100$ )	2.982	0.0480	62.15

Next we consider a similar analysis of these data where the  $\alpha_i$  are no longer fixed but are considered to be random; to emphasize this distinction, the subject-specific effects are now denoted by  $b_i$ ,

$$E(Y_{ij}|b_i) = \beta_1 + \beta_2 \text{Age}_{ij} + b_i,$$

where  $b_i \sim N(0, \sigma^2_b)$ . The results from fitting the random effects model are reported in [Table 9.2](#) and are remarkably similar to those for the fixed effects model. The estimated coefficient for age, 0.0298, indicates that each one year increase in age is associated with a 3.0% relative increase in FEV<sub>1</sub> (adjusted for height<sup>2</sup>). The congruence between the estimates of the effect of age from the fixed and random effects models suggests that any omitted time-invariant characteristics of the children are likely uncorrelated with age and/or that the random effects estimate is primarily based on the within-

subject variation. This congruence can be explained by considering the relative magnitudes of the between- and within-subject sources of variation in these data. For age, 17.2% of the variation is between-subject variation and  $\hat{\rho}_{xx} = 0.172$ . For the response, 69% of the variation is between-subject variation, with  $\hat{\rho}_y = \frac{0.931}{0.931+0.419} \cdot 419 = 0.69$ . Recall from [equation \(9.3\)](#) that the random effects approach is based on an optimal weighted combination of the within-subject and between-subject sources of variation in the response and covariate. As was discussed in Section 9.3, the fixed and random effects approaches yield similar estimates when the weight,

**Table 9.2** Regression (and variance) estimates and standard errors for the random effects model for the  $\log(\text{FEV}_1/\text{height}^2)$  data from the Six Cities Study.

Variable	Estimate	SE	Z
Intercept	-0.355	0.0082	-43.42
Age ( $\times 100$ )	2.981	0.0473	62.99
$\sigma_b^2 (\times 100)$	0.931	0.0849	
$\sigma_e^2 (\times 100)$	0.419	0.0144	

$$w = \frac{(1 - \rho_y)\rho_x}{(1 - \rho_y) + n\rho_y(1 - \rho_x)},$$

is small and close to zero. For the lung function growth data,

$$w \approx \frac{(1 - \rho_y)\rho_x}{(1 - \rho_y) + \bar{n}\rho_y(1 - \rho_x)} = \frac{(1 - 0.69) \times 0.172}{(1 - 0.69) + 6.67 \times 0.69 \times (1 - 0.172)} = 0.013,$$

where  $\bar{n} = \frac{1}{N} \sum_{i=1}^N n_i = 6.67$ . This implies that the estimated effect of age from the random effects model gives approximately 99% of the weight to the fixed effects estimate; that is, estimation of the effect of age is based almost entirely on the within-subject variation in the response and covariate. Although there is some discernible between-subject variation in age (approximately 1/5 of the variation) in these data, its importance for estimation of the effect of age is almost completely downweighted by the magnitude of  $\rho_y$  (the proportion of between-subject variation in the response). Thus, in this instance, the random effects model is a very close approximation to the fixed effects model. Therefore it should not be too surprising that they yield similar estimates of effects and also very similar standard errors. However, we caution that this similarity in estimates between the fixed and random effects models cannot be expected in general.

Although there is almost perfect congruence between the fixed and random effects models for these data, for illustrative purposes we consider a mixed effects model that allows for separate estimation of the effect of age from the between- and within-subject sources of information. The model is

$$E(Y_{ij}|b_i) = \beta_1 + \beta_2^{(L)}(\text{Age}_{ij} - \bar{\text{Age}}_i) + \beta_2^{(C)}\bar{\text{Age}}_i + b_i,$$

where  $b_i \sim N(0, \sigma^2_b)$ . This analysis yields  $\hat{\beta}_2^{(L)}(x 100) = 2.982$  (SE = 0.0480) and  $\hat{\beta}_2^{(C)}(x 100) = 2.923$  (SE = 0.2901). A formal comparison of these two estimates of the effect of age, testing  $H_0: \beta_2^{(L)} = \beta_2^{(C)}$ , produces a chi-squared statistic of 0.04, with 1 df,  $p > 0.80$ . Thus there is no evidence of any conflict between the longitudinal and cross-sectional information. This provides additional justification for the combined estimate of the effect of age from the mixed effects model analysis reported in [Table 9.2](#). Finally, we note that the mixed effects model analysis can be extended in a natural way to allow for randomly varying intercepts and slopes. This allows for a more flexible model for the marginal covariance among the repeated measures of the response (see Section 8.3). Although the addition of randomly varying slopes leads to a discernible improvement in the fit of the model to these data, the results of such an analysis yield an estimate of the effect of age (x 100) of 3.010 (SE = 0.0661) that is very similar to that reported in [Table 9.2](#). However, we note that the standard error is approximately 40% larger and may provide a more realistic estimate of the true sampling variability.

# 9.7 COMPUTING: FITTING LINEAR FIXED EFFECTS MODELS USING PROC GLM IN SAS

To fit linear fixed effects models, we can use PROC GLM in SAS. The GLM procedure fits general linear models using ordinary least squares (OLS) or maximum likelihood estimation under the assumption of normal errors. To fit a fixed effects model, dummy or indicator variables for subjects must be included in the model to jointly estimate  $\beta$  and the  $\alpha_i$ . Linear regression with indicator variables for subjects yields the same estimate of  $\beta$  that would be obtained from the OLS regression of the *mean-centered* response and covariates.

In a regression model that includes an intercept term, say  $\beta_1$ ,  $N - 1$  indicator variables must be included (where  $N$  denotes the number of subjects). That is, we include an indicator variable for each subject except the last; the last subject then becomes the “reference.” These indicator variables can be created automatically by including the subject identifier on the CLASS statement in PROC GLM. Alternatively, in a regression model that omits an intercept term,  $N$  indicator variables must be included. These indicator variables can also be created automatically by including the subject identifier on the CLASS statement in PROC GLM. One small advantage of the latter approach is that the  $N$  subject-specific effects (the  $\alpha_i$ 's) are directly estimated. In the former approach, the subject-specific effect for the last subject (the “reference”) is  $\beta_1$ , while the subject-specific effects for the remaining subjects are  $\beta_1 + \alpha_i$  ( $i = 1, \dots, N - 1$ ). There is ordinarily little interest in the subject-specific effects; they are regarded as nuisance parameters in the fixed effects model formulation.

We note that estimation of  $N$  subject-specific effects increases the computations required for fitting the fixed effects model. This may be of concern in studies where  $N$  is very large. Fortunately, it is possible to reduce this computational burden by not explicitly estimating the  $N$  subject-specific effects; instead, only the covariate effects are estimated. Without explaining the technical details, this is achieved by removing the subject identifier from the MODEL and CLASS statements and including the subject identifier on the ABSORB statement instead. This will yield identical estimates (and standard errors) of the covariate effects but greatly reduce the computation time. The only practical difference in the output is that no estimates of the subject-specific effects are produced. As mentioned earlier, this is usually not a concern as they are regarded as nuisance parameters in the fixed effects model.

For example, to fit a fixed effects model to longitudinal data from a study of aging, we can use the illustrative SAS commands given in [Tables 9.3](#) and [9.4](#). Next we present a brief description of the most salient parts of the command syntax used in the illustrations in [Tables 9.3](#) and [9.4](#).

**Table 9.3** Illustrative commands for fixed effects model using PROC GLM in SAS.

---

```
PROC GLM;
  CLASS id;
  MODEL y=age id / SOLUTION;
```

---

**Table 9.4** Illustrative commands for fixed effects model using the ABSORB statement in PROC GLM in SAS.

---

```
PROC GLM;
  ABSORB id;
  MODEL y=age / SOLUTION;
```

---

PROC GLM <options>;

This statement calls the procedure GLM in SAS. PROC GLM is a very versatile procedure for fitting general linear models, including standard linear regression and ANOVA models.

CLASS variables;

The CLASS statement is used to identify all variables that are to be treated as categorical. By default, this statement will create indicator variables for each listed variable using a reference group coding, with the last level (where “last” refers to the level with the largest alphanumeric value) treated as the reference group. In [Table 9.3](#), the subject identifier, id, is listed on the CLASS statement and this automatically yields  $N - 1$  indicator variables (where  $N$  denotes the number of subjects) when the subject identifier, id, is also included as a covariate in the MODEL statement in PROC GLM.

ABSORB variables;

When a variable is included in the ABSORB statement, the effect of that variable is removed before the construction and solution of the remainder of the model. Absorption can provide a large reduction in computer time and memory requirements, especially when the absorbed effect has a large number of levels. In [Table 9.4](#) the subject identifier, id, is listed on the ABSORB statement. This avoids the explicit estimation of the subject-specific effects,  $\alpha_i$ , in the fixed effects model. As a result it is no longer necessary to include the subject identifier, id, in the MODEL statement. Note: The data set must be sorted by the variable(s) in the ABSORB statement.

MODEL response = <effects> / <options>;

The MODEL statement specifies the response variable and the covariate effects. The covariate effects can include both discrete (defined in the CLASS statement) and quantitative (excluded from the CLASS statement) covariates. By default, PROC GLM includes a column of 1's for the intercept in the model. In [Table 9.3](#) the MODEL statement includes age as a time-varying quantitative (excluded from the CLASS statement) covariate and id as a discrete (defined in the CLASS statement) covariate or factor. Because id is included in the CLASS statement in [Table 9.3](#), this automatically yields  $N - 1$  indicator variables (where  $N$  denotes the number of subjects) for the subject-specific effects. In [Table 9.4](#) the MODEL statement only includes age as a time-varying quantitative (excluded from the CLASS statement) covariate. Because id is included in the ABSORB statement in [Table 9.4](#), explicit estimation of the subject-specific effects,  $\alpha_i$ , is avoided.

# 9.8 COMPUTING: DECOMPOSITION OF BETWEEN-SUBJECT AND WITHIN-SUBJECT EFFECTS USING PROC MIXED IN SAS

To fit the linear mixed effects model described in Section 9.4, we can use PROC MIXED in SAS. To allow for simultaneous estimation of both the cross-sectional and longitudinal effects of a time-varying covariate, we simply include both the mean of the covariate (averaged over time) and the mean-centered covariate in the analysis. The estimated coefficient for the mean of the covariate yields an estimate of the cross-sectional effect. The estimated coefficient for the mean-centered covariate yields an estimate of the longitudinal effect. Fitting linear mixed effects models that decompose between-subject and within-subject effects is only slightly more complicated than fitting standard linear mixed effects models.

For example, to fit a mixed effects model to longitudinal data from a study of aging, we can use the illustrative SAS commands given in [Table 9.5](#). In this illustration we estimate both the cross-sectional and longitudinal effects of aging and also the effect of gender, a time-invariant covariate. Recall that the fixed effects model precludes estimation of the effects of any time-invariant covariates.

**Table 9.5** Illustrative commands for decomposition of between- and within-subject effects using PROC MIXED in SAS.

---

```
PROC MEANS DATA=one NWAY;
  CLASS id;
  VAR age;
  OUTPUT OUT=two MEAN=mage;
PROC SORT DATA=one;
  BY id;
PROC SORT DATA=two;
  BY id;
DATA three;
  MERGE one two;
  BY id;
  cage=age-mage;
RUN;

PROC MIXED DATA=three;
  CLASS id gender;
  MODEL y=gender cage mage / SOLUTION;
  RANDOM INTERCEPT / SUBJECT=id;
```

---

The first part of the command syntax used in the illustration in [Table 9.5](#) is for the calculation of the mean of the time-varying covariate, age. The PROC MEANS procedure in SAS (with the NWAY option) can be used to calculate the mean age for each subject, where averaging is over repeated observations within a subject. Using an OUTPUT statement, the mean age, denoted `mage` in [Table 9.5](#), is written to a second SAS data-set (two). The original SAS data-set (one) and the second SAS data-set (two) are then sorted according to the subject-identifier, `id`. Finally, a third SAS data-set (three) is created by merging these two data-sets. Using the third SAS data-set (three) the mean-centered age variable, denoted `cage` in [Table 9.5](#),

`cage = age - mage,`  
is calculated for each observation. The mean age and mean-centered age variables are then used as

covariates in the mixed effects model analysis.

The second part of the command syntax used in the illustration in [Table 9.5](#) is for the fitting of a random effects (random intercepts only) model that includes the effect of gender and both cross-sectional and longitudinal effects of age. Using a CONTRAST statement, the cross-sectional and longitudinal estimates of the effects of aging can be compared.

## 9.9 FURTHER READING

Allison (2005) provides a remarkably clear and well-organized guide to fixed effects models for longitudinal data, with numerous examples of how to implement the methods using standard statistical software. Ware et al. (1990) present an accessible discussion of regression models for longitudinal data that incorporate separate parameters for cross-sectional and longitudinal effects of aging; also see discussion of discrepancies between longitudinal and cross-sectional effects in Louis et al. (1986). Neuhaus and Kalbfleisch (1998) and Neuhaus (2001) discuss similar issues in the broader context of cluster-correlated data; also see Berlin et al. (1999) and Begg and Parides (2003) for a discussion of the decomposition of between- and within-cluster effects in the analysis of cluster-correlated data.

# Bibliographic Notes

The derivation of the “Hausman test” is described in detail in Hausman (1978). Because a very similar approach was first proposed by Durbin (1954), and separately by Wu (1973), tests based on the comparison of two sets of parameter estimates are also referred to as “Durbin–Wu–Hausman” tests.

The statistical basis for the weights in [equations \(9.3\)](#) and [\(9.4\)](#) is discussed in Scott and Holt (1982); a closely related expression for the weights applied to the between-subject and within-subject estimators can also be found in an influential paper in the econometrics literature by Maddala (1971).

## Problems

**9.1** In a U.S. Centers for Disease Control (CDC) study, conducted in the state of Georgia from 1980 to 1992, linked data on live births to the same mother were obtained (Adams et al., 1997). We focus on a subset of data restricted to 878 mothers for whom five births were identified; these data are from Neuhaus and Kalbfleisch (1998) and are reported in Pan (2002) and Rabe-Hesketh and Skrondal (2008). The main objective of the analysis is to determine the effect of maternal age (the mother’s age at each birth) on infant birth weight. Note that because each mother had five infant births, but at different maternal ages, maternal age is a within-subject or time-varying covariate.

The raw data are stored in an external file: birthwt.dat

Each row of the data set contains the following five variables:

MID Order Wt Age CID

*Note:* The outcome variable Wt records the infant birth weight measured in grams. The variable Order denotes the infant birth order (coded 1–5). The variable Age is the mother’s age (in years) at each of the five recorded births. Finally, the variables MID and CID denote the mother and child identifiers (IDs).

**9.1.1** Let  $Y_{ij}$  denote the birth weight (in grams) of the  $j^{th}$  infant from the  $i^{th}$  mother and  $\text{Age}_{ij}$  denote the  $i^{th}$  mother’s age at the time of the birth of her  $j^{th}$  infant (for  $i = 1, \dots, 878; j = 1, \dots, 5$ ). Fit the following standard mixed effects model with linear trend for maternal age and randomly varying intercepts,

$$Y_{ij} = \beta_1 + \beta_2 \text{Age}_{ij} + b_i + \epsilon_{ij},$$

where  $b_i \sim N(0, \sigma_b^2)$  and  $\epsilon_{ij} \sim N(0, \sigma_\epsilon^2)$ .

**9.1.2** From the results of the analysis for Problem 9.1.1, is there evidence of an association between infant birth weight and maternal age? Present results that support your conclusion.

**9.1.3** What is the interpretation of  $\hat{\beta}_2$ ?

**9.1.4** For a fixed effects analysis of these data, fit the following fixed effects model with linear trend for maternal age:

$$Y_{ij} = \beta_1 + \beta_2^{(FE)} \text{Age}_{ij} + \alpha_i + \epsilon_{ij},$$

where the  $\alpha_i$  are considered fixed, not random effects, and  $\epsilon_{ij} \sim N(0, \sigma_\epsilon^2)$ .

**9.1.5** From the results of the analysis for Problem 9.1.4, is there evidence of an association between infant birth weight and maternal age? Present results that support your conclusion.

**9.1.6** What is the interpretation of  $\hat{\beta}_2^{(FE)}$ ?

**9.1.7** Compare and contrast  $\hat{\beta}_2$  from the mixed effects model in Problem 9.1.1 and  $\hat{\beta}$  from the fixed effects model in Problem 9.1.4. Explain why they might differ.

**9.1.8** Consider a mixed effects model analysis of these data that allows for separate estimation of the cross-sectional and longitudinal effects of maternal age on infant birth weight. Fit the following mixed effects model with separate linear trends for the *cross-sectional* (C) and *longitudinal* (L) effects of maternal age:

$$Y_{ij} = \beta_1 + \beta_2^{(C)} \overline{\text{Age}}_i + \beta_2^{(L)} (\text{Age}_{ij} - \overline{\text{Age}}_i) + b_i + \epsilon_{ij},$$

where  $b_i \sim N(0, \sigma_b^2)$  and  $\epsilon_{ij} \sim N(0, \sigma_\epsilon^2)$ .

**9.1.9** What is the interpretation of  $\hat{\beta}_2^{(C)}$ ?

**9.1.10** What is the interpretation of  $\hat{\beta}_2^{(L)}$ ?

**9.1.11** Compare and contrast  $\hat{\beta}_2^{(C)}$  and  $\hat{\beta}_2^{(L)}$ . Explain why they might differ.

**9.1.12** Construct a formal test of  $H_0: \beta_2^{(C)} = \beta_2^{(L)}$ ; this test can be based on either a Wald test or a likelihood-ratio test formed by comparing nested models. What do you conclude?

**9.1.13** Compare and contrast  $\hat{\beta}_2^{(C)}$  and  $\hat{\beta}_2^{(L)}$  with  $\hat{\beta}_2$  from the mixed effects model in Problem 9.1.1 and  $\hat{\beta}_2^{(FE)}$  from the fixed effects model in Problem 9.1.4.

<sup>1</sup> In the econometrics literature, this mean-centered transformation is often referred to as a “demeaning” transformation; we avoid the use of this term lest we create the impression there is something improper about the subtraction of a mean.

<sup>2</sup> Note that there are several transformations that eliminate  $\alpha_i$  from the model, such as subtraction of the baseline values,  $Y_{ij\cdot}^* = (Y_{ij} - Y_{i1})$  and  $x_{ijk}^* = (X_{ijk} - X_{i1k})$  (see Section 8.8); however, in general, mean-centering is the preferred transformation.

<sup>†</sup> The remainder of this section can be skipped on a first reading. However, the reader who returns to it will find that it yields insights about when the fixed and mixed effects estimators are likely to be similar.

# *Chapter 10*

## *Residual Analyses and Diagnostics*

### **10.1 INTRODUCTION**

The analysis of longitudinal data is not complete without an examination of the residuals. Residuals can be used to assess the adequacy of the fitted model and can also indicate the presence of outliers. Methods for residual analyses are well developed for standard regression settings with independent observations on a univariate response. In principle, many of the same properties of residual analysis can be extended to the longitudinal setting.

## 10.2 RESIDUALS

With longitudinal data we can define a vector of residuals for each individual,

$$(10.1) \quad r_i = Y_i - X_i \hat{\beta}.$$

The vector of residuals has mean zero and provides an estimate of the vector of errors,

$$e_i = Y_i - X_i \beta.$$

The residuals defined in (10.1) can be used to check for any systematic departures from the model for the mean response; they can also form the basis of an assessment of the adequacy of the model for the covariance. For example, a scatterplot of the residuals,

$$r_{ij} = Y_{ij} - X'_{ij} \hat{\beta},$$

against the predicted mean response,

$$\hat{\mu}_{ij} = X'_{ij} \hat{\beta},$$

can be examined for the appearance of any systematic trend. The fitting of a smooth curve (e.g., a *lowess* curve; see Section 3.3) to the scatterplot can often help in judging whether curvature is present. In a correctly specified model, the scatterplot should display no systematic pattern, with a more or less random scatter around a constant mean of zero. Similarly scatterplots of the residuals against selected covariates from the mean model can be examined for any systematic trends. Such a trend may indicate the omission of a quadratic term or the need for transformation of the covariate.

For most practical purposes, graphical displays of the residuals can be used to detect discrepancies in the model for the mean response or the presence of outlying observations that require further investigation. However, there are two properties of the residuals from an analysis of longitudinal data that must be kept in mind. First, the components of the vector of residuals,

$$r_i = Y_i - X_i \hat{\beta},$$

are correlated and do not necessarily have constant variance. Recall that the mean of the residuals is zero, mimicking the mean of the vector of errors,

$$e_i = Y_i - X_i \beta.$$

In contrast, the covariance of the residuals is not identical to the covariance of the errors. However, for all practical purposes we can approximate the covariance of the residuals by

$$\text{Cov}(r_i) \approx \text{Cov}(e_i) = \Sigma_i.$$

Because the residuals have approximate covariance matrix,  $\Sigma_i$ , this has important implications for the examination of plots of the residuals. First, because the variance is not necessarily constant, the scatterplot of the residuals against the predicted values, or against time, will not necessarily have a constant range. As a result standard residual diagnostics for examining either the homogeneity of the residual variance or autocorrelation among the residuals should be avoided altogether. Second, although residuals from a univariate linear regression are uncorrelated with the covariates, the residuals from a regression analysis of longitudinal data may be correlated with the covariates. As a result there may be an apparent systematic trend in the scatterplot of the residuals against a selected covariate.

## 10.3 TRANSFORMED RESIDUALS

To circumvent some of the aforementioned problems with the use of residuals from longitudinal data based on (10.1), we can transform the residuals. There are many possible ways to transform the residuals. In general, it would be desirable to standardize and, for lack of a better term, “de-correlate” the residuals so that they mimic residuals from a standard linear regression. That is, we would like to transform the residuals so that they have constant variance and zero correlation. This can be achieved using a simple and well-known method called the *Cholesky decomposition* (or *Cholesky factorization*).

Given an estimate of the approximate covariance matrix for the residuals,  $\hat{\Sigma}_i$ , the Cholesky decomposition of  $\hat{\Sigma}_i$  can be used to create a lower triangular matrix,  $L_i$ , such that

$$\hat{\Sigma}_i = L_i L_i';$$

note that a lower triangular matrix is simply one with all zeros above the diagonal. We can then use the matrix  $L_i$  or, more specifically,  $L_i^{-1}$  to take us from a set of correlated residuals with heterogeneous variances to a set of transformed residuals,

$$(10.2) \quad r_i^* = L_i^{-1} r_i = L_i^{-1} (Y_i - X_i \hat{\beta}),$$

which are uncorrelated and have unit variance.

Interestingly, the transformation used in (10.2) leads to a set of transformed residuals with appealing interpretations in the longitudinal setting. For example, the first element of

$$r_i^* = L_i^{-1} (Y_i - X_i \hat{\beta})$$

is the standardized residual for the first repeated measurement (often the baseline observation). In contrast, the second through last transformed residuals represent standardized deviations from the conditional mean of the response given all previous observations. For example, the  $k^{th}$  transformed residual is an estimate of

$$\frac{Y_{ik} - E(Y_{ik}|Y_{i1}, \dots, Y_{ik-1})}{\sqrt{\text{Var}(Y_{ik}|Y_{i1}, \dots, Y_{ik-1})}}.$$

This property of the transformed residuals for longitudinal data is not obvious; it requires a deeper understanding of matrix algebra than is assumed throughout this book.

Given the set of transformed residuals,  $r_i^*$ , all of the usual residual diagnostics for standard linear regression can be applied. For example, we can construct a scatterplot of the transformed residuals,  $r_{ij}^*$ , versus the transformed predicted values,  $\hat{\mu}_{ij}^*$ , where

$$\hat{\mu}_i^* = L_i^{-1} \hat{\mu}_i = L_i^{-1} X_i \hat{\beta}.$$

In a correctly specified model, this scatterplot should display no systematic pattern, with a random scatter around a constant mean of zero and with a constant range for varying  $\hat{\mu}_{ij}^*$ . Similarly we can construct a scatterplot of the transformed residuals versus selected transformed covariates. With longitudinal data, a scatterplot of the transformed residuals versus transformed time (or age) can be particularly useful for assessing the adequacy of the model assumptions about patterns of change in the mean response over time. Finally, the transformed residuals also make it somewhat easier to identify skewness and potential outliers that require further investigation. A normal quantile plot (or so-called quantile-quantile or Q-Q plot) of the transformed residuals can be used to assess the normal distribution assumption and to identify outliers. That is, on the basis of the ranks of the transformed residuals, we can plot the sample quantiles of the residuals against the quantiles expected if they have a normal distribution. If the residuals depart discernibly from a straight line, then the assumption of normality may not be tenable. For example, skewness is usually indicated by a bow-shaped pattern in the normal quantile plot; outliers will appear as “stragglers,” far from the ends of the line.

When statistical software is available that automates the production of residual diagnostics for standard linear regression, the following procedure can be used. Given the estimated covariance,  $\hat{\Sigma}_i$ ,  $L_i$  can be obtained from the Cholesky decomposition of  $\hat{\Sigma}_i$ . Then a transformed response vector and covariate matrix can be constructed as follows:

$$Y_i^* = L_i^{-1} Y_i; \quad X_i^* = L_i^{-1} X_i.$$

Finally, the generalized least squares (GLS) estimate of  $\beta$ , from the regression of  $Y_i$  on  $X_i$  (with estimated covariance matrix,  $\hat{\Sigma}_i$ ), can be re-estimated from the ordinary least squares (OLS) regression of  $Y_i^*$  on  $X_i^*$ . That is, any standard linear regression program can be used to model the dependence of  $Y_i^*$  on  $X_i^*$ , and all the built-in residual diagnostics from the resulting OLS regression can be examined to check the adequacy of the model. Thus, once a model has been selected for the mean and the covariance, standard regression diagnostics for independent observations (with homogeneous variance) can be applied by re-fitting a standard linear regression of  $Y_i^*$  on  $X_i^*$  and making use of available residual diagnostic tools.

As mentioned earlier, the transformed residuals are useful for detecting outlying *observations*. They can also be used to detect outlying *individuals*. For each individual we can calculate a summary measure of multivariate distance between their observed and fitted responses, based on the *Mahalanobis distance*,

$$(10.3) \quad d_i = r_i^{*'} r_i^*.$$

If the model is correctly specified, the distances given by (10.3) have an approximate chi-squared distribution with degrees of freedom (df) equal to the dimension of  $r_i^*$  (i.e.,  $df = n_i$ , the number of repeated measurements on the  $i^{th}$  subject). Outlying individuals will have distances,  $d_i$ , that have small associated  $p$ -values. The  $p$ -values provide a common metric for comparing and detecting large values of  $d_i$ , corresponding to unusual or outlying individuals, when the number of repeated measurements varies across subjects. We should add a word of caution concerning the interpretation of these  $p$ -values. Because the major focus is on the most extreme values of  $d_i$ , the distribution of these extremes (e.g., the distribution of the maximum  $d_i$ ) is somewhat more complicated than a chi-squared distribution with  $n_i$  degrees of freedom. In principle, a Bonferroni correction to the  $p$ -values could be applied (e.g., multiplying the  $p$ -values by the sample size,  $N$ ); however, the Bonferroni correction is known to be very conservative. In general, we recommend that the  $p$ -values be used as a common metric for comparing  $d_i$  when the number of repeated measurements varies across subjects, while recognizing that  $p$ -values less than 0.05 occur with predictable regularity (when the sample size is 200, the expected number is  $200 \times 0.05 = 10$ ).

So far much of the discussion of residual diagnostics has focused on the adequacy of the model for the mean response. As alluded to above, the adequacy of the variance assumption can be informally assessed by examining the scatterplot of the transformed residuals versus the transformed predicted values and/or time. In a correctly specified model for the variance, the range of the transformed residuals should be approximately constant over (transformed) time and for varying  $\hat{\mu}_{ij}^*$ . A more informative plot is obtaining by considering the scatterplot of the absolute values of the transformed residuals,  $|r_{ij}^*|$ , versus (transformed) time and/or  $\hat{\mu}_{ij}^*$ . If the assumed model for the variance is adequate, there should be no systematic trend. The fitting of a smooth curve to the scatterplot (e.g., a lowess curve) can often help in judging whether any curvature is present. The fitted curve should display no systematic departures from a horizontal line centered at approximately 0.8; note that, if the transformed residuals are assumed to be normal, with mean zero and unit variance, then the mean of the absolute values of the residuals is 0.798. Finally, an informal check on the overall adequacy of the model for the covariance, both the models for the variances and correlations, is provided by a smoothed plot of the so-called empirical semi-variogram. The definition of the empirical semi-variogram, and a description of its use as a diagnostic tool, are given in Section 10.5.

## 10.4 AGGREGATING RESIDUALS

An acknowledged difficulty with conventional residual diagnostics is that they are somewhat subjective in nature. What appears to be a random scatter to one individual, might be considered evidence of systematic trend to another. It can be difficult to discern whether an apparent trend in a scatterplot of the residuals reflects some aspect of model misspecification or is simply a reflection of natural variation. McCullagh and Nelder (1989, pp. 392–393) aptly summarize this problem when they state that “the practical problem is that any finite set of residuals can be made to yield some kind of pattern if we look hard enough, so that we have to guard against over-interpretation.”

Recently, model-checking techniques based on “cumulative sums” and “moving sums” of residuals have been developed to help discern the “signal” from the “noise” in scatterplots of residuals. The basic idea is to aggregate or sum the residuals over certain coordinates. The coordinates typically used for these sums are the individual covariates (e.g.,  $X_{ijk}$ , the  $k^{\text{th}}$  covariate) and the fitted values,  $X'_{ij}\hat{\beta}$ . The chief advantage of working with sums of residuals is that a reference distribution is available to ascertain their natural variation. That is, we can compare the *observed* sum of the residuals, both graphically and numerically, to a reference distribution under the assumption of a correctly specified model for the mean response. This allows us to determine whether any apparent pattern in the observed sum of the residuals is evidence of a systematic trend or simply due to natural variation. This way we can remove a large degree of subjectivity from the assessment of graphical displays of residuals, placing residual diagnostics on a more objective footing.

Recall from Section 10.2 that the raw residuals are defined as the difference between the observed and fitted values of the responses,

$$r_{ij} = Y_{ij} - X'_{ij}\hat{\beta}.$$

If the model for the mean is correctly specified, these residuals are centered at zero. To check the functional form of any covariate, say  $X_{ijk}$ , the  $k^{\text{th}}$  covariate, we can define the cumulative sum of the residuals over values of  $X_{ijk}$ ,

$$W_k(x) = \frac{1}{\sqrt{N}} \sum_{i=1}^N \sum_{j=1}^{n_i} I(X_{ijk} \leq x) r_{ij},$$

where  $I(\cdot)$  is the indicator function. For any given value  $x$ ,  $W_k(x)$  is the sum of residuals for all values of  $X_{ijk}$  less than or equal to  $x$ . The process  $W_k(x)$  is a step function with possible jumps, either increases or decreases, at all of the distinct values for  $X_{ijk}$ . In addition we can construct the cumulative sum of residuals over the fitted values, denoted by  $W_f(x)$ ,

$$W_f(x) = \frac{1}{\sqrt{N}} \sum_{i=1}^N \sum_{j=1}^{n_i} I(X'_{ij}\hat{\beta} \leq x) r_{ij}.$$

The cumulative sum,  $W_k(x)$ , can then be used to provide both a graphical and numerical assessment of the functional form of the covariate. For example, we can construct a plot of  $W_k(x)$  versus  $x$ , where for any value of  $x$  on the horizontal axis, the corresponding value of  $W_k(x)$  on the vertical axis is the cumulative sum of the residuals for covariate values of  $X_{ijk}$  less than or equal to  $x$ . Evidence of systematic trend in this plot suggests that the functional form of the covariate (e.g., linearity) is not correctly specified and may indicate that a transformation of the covariate or the inclusion of polynomials is required. The cumulative sum,  $W_f(x)$ , is useful for assessing the assumption of linearity or, more generally, the functional form of the relationship between the response and the covariates. Any evidence of systematic trend in the latter plot might suggest that either a transformation of  $Y$  (i.e., a transformation of the response) or of  $E(Y|X)$  (i.e., a transformation of the conditional mean of  $Y$ ) is necessary.

How do we determine if there is any evidence of systematic trend in a plot of the cumulative sum of residuals? If the assumed model for the mean response has been correctly specified, the cumulative residual process should be centered at zero and behave like a zero-mean Gaussian (or normal) process. This zero-mean Gaussian process provides a reference for ascertaining whether

any pattern in the observed cumulative residual process is systematic or simply due to the natural variation of the process. Therefore an assessment of model adequacy can be based on comparing the pattern of the *observed* cumulative residual process with computer simulated realizations from the *expected* process, a zero-mean Gaussian process (the null distribution). It is relatively straightforward to simulate realizations from the null distribution of the cumulative sum of residuals, although the technical details are omitted here; the interested reader is referred to the article by Lin, Wei, and Ying (2002).

Thus, in practical terms, the null distributions of  $W_k(x)$  and  $W_f(x)$  can be approximated through computer simulation of the zero-mean Gaussian processes, denoted by  $\widehat{W}_k(x)$  and  $\widehat{W}_f(x)$ , respectively. Then, to assess whether any apparent trend in the *observed* cumulative sum of residuals reflects systematic trend rather than chance fluctuations, we can superimpose a number of realizations from the appropriate Gaussian process. To the extent that the curves generated from the null distribution tend to be closer to and intersect zero more often than the observed curve, this provides evidence of lack of fit. This assessment can be put on a more formal footing by comparing the maximum absolute value of the observed cumulative sum to the maximum absolute value from a large number of realizations (e.g., 10,000) from the null distribution. For example, by comparing  $\max|W_f(x)|$ , the maximum absolute value of the observed cumulative sum, to  $\max|\widehat{W}_f(x)|$  for each realization from the null distribution, a  $p$ -value can be constructed based on the proportion of times that  $|\widehat{W}_f(x)| \geq \max|W_f(x)|$ . The latter is referred to as a “supremum” test and provides an omnibus test of model adequacy with respect to the relevant coordinate (e.g., the fitted values or a particular covariate). If the  $p$ -value is very small (e.g., less than 0.05 or 0.01), this indicates that the observed cumulative residual process is extreme, suggesting that the fit of the model to the mean response can be improved.

There is an alternative way to aggregate the residuals by using a “moving sum” rather than a “cumulative sum.” We can define a moving sum of residuals, with “window” of width  $b$ , as follows:

$$W_k(x, b) = \frac{1}{\sqrt{N}} \sum_{i=1}^N \sum_{j=1}^{n_i} I(x - b \leq X_{ijk} \leq x) r_{ij}.$$

This represents the sum of residuals in blocks of window size,  $b$ . Similarly, to assess linearity (or the functional form of the relationship between the response and covariates), we can define a moving sum of residuals with respect to the fitted values,

$$W_f(x, b) = \frac{1}{\sqrt{N}} \sum_{i=1}^N \sum_{j=1}^{n_i} I(x - b \leq X'_{ij}\hat{\beta} \leq x) r_{ij}.$$

A potential advantage of using a moving sum of residuals is that the process is less influenced by the residuals associated with small values of the covariate. One disadvantage of moving sums, however, is that they require a somewhat arbitrary choice of window size,  $b$ . Based on results from simulation studies, Lin, Wei, and Ying (2002) suggest that the optimal choice of  $b$  is approximately the range of the lower half of the covariate values.

Finally, an appealing property of these graphical and numerical methods based on cumulative and moving sums of residuals is that they are valid regardless of the true joint distribution of the longitudinal response vector. In particular, these model-checking techniques do not require correct specification of the covariance among the responses. As such, these graphical and numerical techniques for assessing the model for the mean response are relatively robust to assumptions about the distribution of the responses and assumptions about the covariance among the repeated measures.

## 10.5 SEMI-VARIOGRAM

*t<sub>1</sub>*  
*t<sub>2</sub>*?  
*t<sub>3</sub>*?

Historically the semi-variogram has been widely used in spatial statistics to represent the covariance structure in geostatistical data. Unlike two-dimensional spatial data, the coordinates for longitudinal data are along a single dimension, namely time. For longitudinal data the semi-variogram is defined as one-half the expected squared difference between residuals obtained on the same individual. The semi-variogram, denoted by  $\gamma(h_{ijk})$ , is given by

$$(10.4) \quad \gamma(h_{ijk}) = \frac{1}{2} E(r_{ij} - r_{ik})^2,$$

where  $h_{ijk}$  is the time elapsed between the  $j^{th}$  and  $k^{th}$  repeated measurement on the  $i^{th}$  individual. Since the residuals have mean zero, the semi-variogram given by (10.4) can be expressed as

$$\begin{aligned} \gamma(h_{ijk}) &= \frac{1}{2} E(r_{ij} - r_{ik})^2 \\ &= \frac{1}{2} E(r_{ij}^2 + r_{ik}^2 - 2r_{ij}r_{ik}) \\ &= \frac{1}{2} \text{Var}(r_{ij}) + \frac{1}{2} \text{Var}(r_{ik}) - \text{Cov}(r_{ij}, r_{ik}). \end{aligned}$$

Although the semi-variogram can be used to suggest appropriate models for the covariance, here we simply use it as a diagnostic tool for assessing the adequacy of a selected model for the covariance. When the semi-variogram is applied to the transformed residuals,  $r_{ij}^*$ , it simplifies to

$$\gamma(h_{ijk}) = \frac{1}{2} \text{Var}(r_{ij}^*) + \frac{1}{2} \text{Var}(r_{ik}^*) - \text{Cov}(r_{ij}^*, r_{ik}^*) = \frac{1}{2}(1) + \frac{1}{2}(1) - 0 = 1.$$

Thus, in a correctly specified model for the covariance, the plot of the semi-variogram for the transformed residuals versus the time elapsed between the corresponding observations should fluctuate randomly around a horizontal line centered at 1.

The empirical or sample semi-variogram,  $\gamma(h)$ , is simply defined as one-half the average squared difference between pairs of residuals on the same individual whose corresponding observations are  $h$  units apart (and where the average is taken over all pairs of observations for which  $h_{ijk} = h$ ). With inherently unbalanced longitudinal data, where subjects are not all measured at the same set of occasions, the empirical semi-variogram can be estimated by fitting a smooth curve (e.g., a lowess curve; see Section 3.3) to the scatterplot of the observed half squared differences between residuals obtained on the same individual and the corresponding time lags. In a correctly specified model for the covariance matrix, a smooth plot of the empirical semi-variogram for the transformed residuals should be centered at 1 and display no systematic curvature. However, the construction of a smooth plot of the empirical semi-variogram requires extra care. Because the empirical semi-variogram is based on the squared differences between pairs of residuals it can be very sensitive to outliers; furthermore, because each residual contributes to  $n_i - 1$  squared differences between pairs of residuals on the same individual, a single outlier can have an inordinate influence at several different time lags.

Finally, recall that in the linear mixed effects model  $\Sigma_i$  has a characteristic random effects structure given by

$$\Sigma_i = Z_i G Z_i' + \sigma^2 I_{n_i}.$$

Transformed residuals from the linear mixed effects model can be obtained by taking the Cholesky decomposition of  $\Sigma_i = Z_i \hat{G} Z_i' + \hat{\sigma}^2 I_{n_i}$ . The adequacy of the random effects covariance structure can be assessed from the plot of the empirical semi-variogram for the transformed residuals. In addition, with linear mixed effects models we can obtain predictions of the random effects (empirical BLUPs) and examine their distribution for any evidence of extremes or outliers, perhaps representing individuals whose subject-specific response profiles are somewhat unusual. Because the empirical BLUPs are known to be heavily influenced by the normal distribution assumption for the random effects, we caution that histograms and normal quantile plots of the empirical BLUPs should not be used to assess the adequacy of the normal distribution assumption for the random effects. Furthermore, because the empirical BLUPs have been “shrunk” toward the population fixed effects,  $\beta$ , their distribution does not accurately represent the distribution of the random effects (e.g., due to “shrinkage” toward the population mean, empirical BLUPs have smaller variance).

## 10.6 CASE STUDY

In Section 8.8 we presented the results of analyses of body fat accretion from a prospective study of the development of obesity in a cohort of girls. In this section we examine the residuals from the fitted model to assess the overall adequacy of the model. The residual analyses are also used to detect individuals with unusual response profiles.

# Study of Influence of Menarche on Changes in Body Fat Accretion

Recall that the data are from a prospective longitudinal study examining changes in body fat before and after menarche in a cohort of 162 girls from the MIT Growth and Development Study (Bandini et al., 2002; Phillips et al., 2003). At the start of the study, all of the girls were pre-menarcheal and non-obese. They were followed over time according to a schedule of annual measurements until four years after menarche. At each examination, a measure of body fatness, percent body fat (%BF) was derived from three basic measurements of (1) body weight, (2) height, and (3) bioelectric impedance resistance.

In Section 8.8 we presented analyses of the changes in percent body fat before and after menarche. For these analyses “time” was coded as time since menarche and could be positive or negative. We considered the hypothesis that percent body fat increases linearly with age, but with different slopes before and after menarche. Specifically, we assumed that each girl had a piecewise linear spline growth curve with a knot at the time of menarche and fitted the following linear mixed effects model

$$E(Y_{ij}|b_i) = \beta_1 + \beta_2 t_{ij} + \beta_3 (t_{ij})_+ + b_{1i} + b_{2i} t_{ij} + b_{3i} (t_{ij})_+,$$

where  $t_{ij}$  denotes the time of the  $j^{th}$  measurement on the  $i^{th}$  subject before or after menarche (i.e.,  $t_{ij} = 0$  at menarche),  $(t_{ij})_+ = t_{ij}$  if  $t_{ij} > 0$  and  $(t_{ij})_+ = 0$  if  $t_{ij} \leq 0$ . In this model each girl’s growth curve can be described with an intercept and two slopes, one slope for the changes in response before menarche, another slope for the changes in response after menarche.

The REML estimates of the fixed effects are displayed in [Table 10.1](#) (the REML estimates of the variance components for the random effects are displayed in [Table 8.7](#) in Chapter 8). The main goal of the analysis was to assess whether the population slopes for fat accretion differ before and after menarche. Based on the magnitude of the estimate of  $\beta_3$ , relative to its standard error, it was concluded that there was a significant difference between the slopes before and after menarche. In particular, the estimated pre-menarcheal slope is rather shallow (0.42) and indicates that the annual rate of body fat accretion is less than 0.5%. In contrast, the estimated post-menarcheal slope is 2.46 and indicates that the annual rate of body fat accretion is approximately 2.5%, almost six times higher than the corresponding rate in the pre-menarcheal period.

**Table 10.1** Estimated regression coefficients (fixed effects) and standard errors for the piecewise linear model for the percent body fat data.

Variable	Estimate	SE	Z
Intercept	21.3614	0.5646	37.84
Time	0.4171	0.1572	2.65
(Time) <sub>+</sub>	2.0471	0.2280	8.98

Next we use residual diagnostics to assess the adequacy of the fitted model. Based on the Cholesky decomposition of the estimated covariance matrix,  $\hat{\Sigma}_i$ , we can calculate transformed residuals,

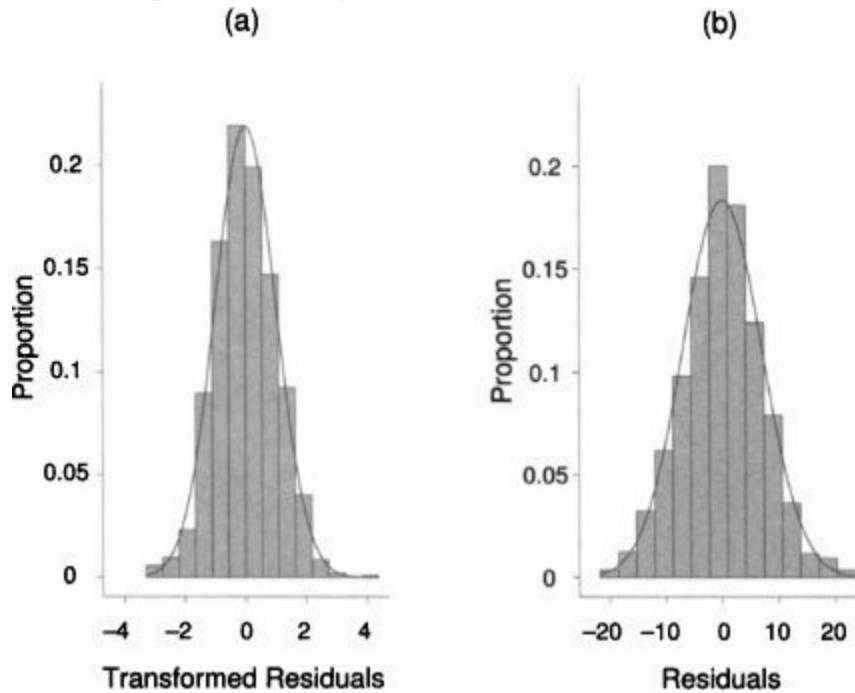
$$r_i^* = L_i^{-1} r_i = L_i^{-1} (Y_i - X_i \hat{\beta}),$$

where  $\hat{\Sigma}_i = L_i L_i'$ . For illustrative purposes we also examine the untransformed residuals and compare the diagnostic plots based on these two types of residuals.

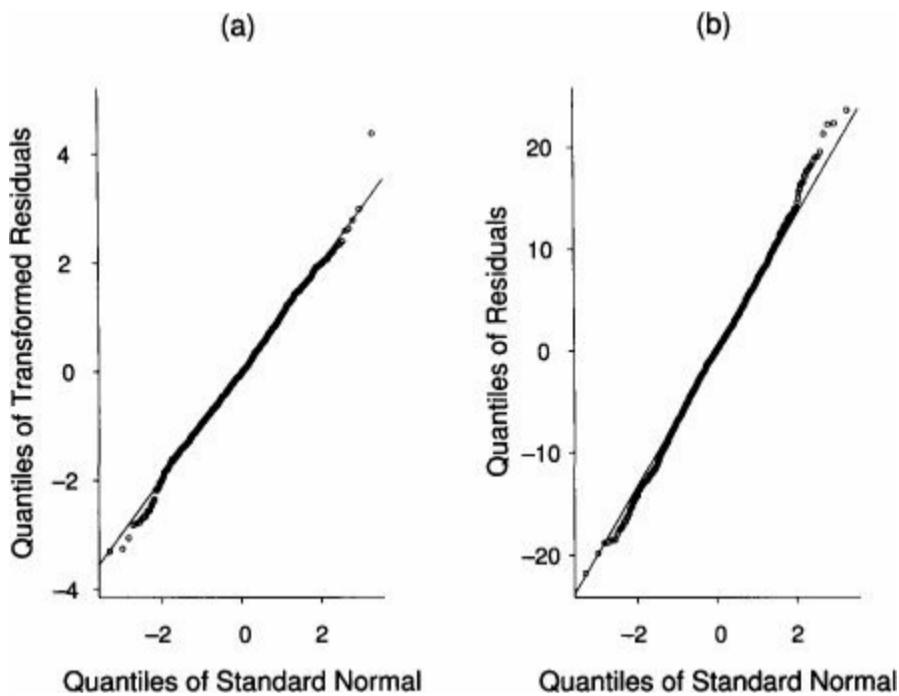
Histograms of the transformed and untransformed residuals are presented in [Figure 10.1](#) and do not indicate any discernible skewness. In addition the normal quantile plots of the residuals do not display any systematic departures from a straight line (see [Figure 10.2](#)). The quantile plot of the transformed residuals does reveal one very extreme observation. This observation corresponds to a measurement on a girl, with subject ID = 128, obtained prior to menarche. Approximately two years prior to menarche, this girl’s percent body fat increased to 44% (from 27% one year earlier); this is an extreme value, over 3 standard deviations above the mean for pre-menarcheal percent body fat. However, the observation is not a recording error and this girl had subsequent measurements of

percent body fat of 40% and 41% at the next two occasions. Overall, the number of extreme residuals highlighted by [Figure 10.2](#) is not more than what we would expect due to chance, given a total of 1049 observations. From an examination of the residuals in [Figures 10.1](#) and [10.2](#), there is no evidence to suggest any discernible skewness and the normal assumption appears to be tenable.

**Fig. 10.1** Histogram, with normal density overlaid, of (a) the transformed residuals, and (b) the untransformed residuals, for the percent body fat data.



**Fig. 10.2** Normal quantile plot of (a) the transformed residuals, and (b) the untransformed residuals, for the percent body fat data.



Before considering scatterplots of the transformed (and untransformed) residuals, we illustrate how the transformed residuals can be used to identify unusual individuals. We can calculate the Mahalanobis distance,

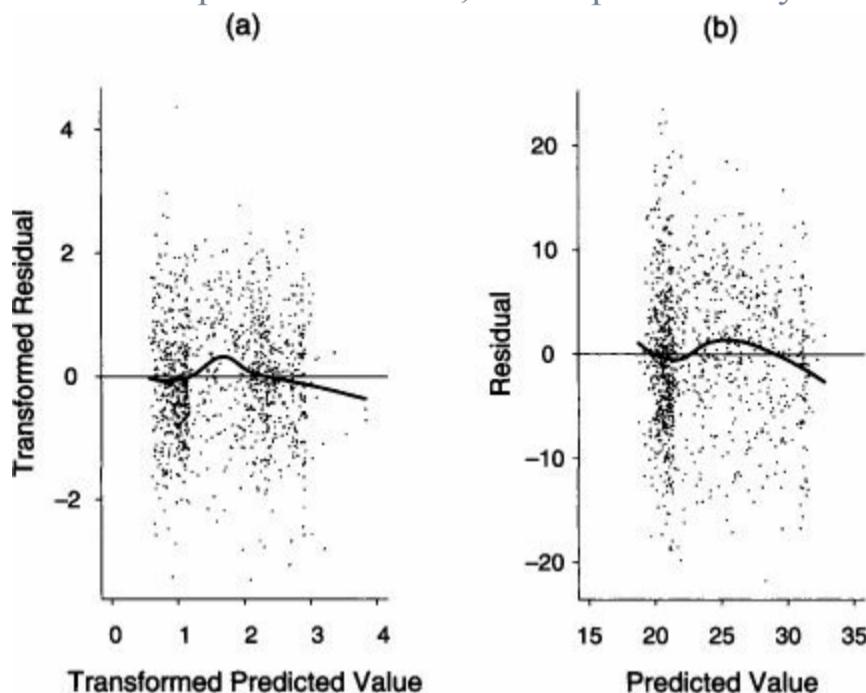
$$d_i = r_i^{*'} r_i^*$$

for each girl and then compare the values to reference chi-squared distributions with degrees of freedom (df) equal to the dimension of  $r_i^*$  (i.e.,  $df = n_i$ , the number of repeated measurements obtained on each girl). For each girl we calculated  $d_i$  and its associated  $p$ -value. There were 7 girls whose  $d_i$  yielded  $p$ -values less than 0.05 and 2 girls with  $p$ -values less than 0.01. Given that the sample is composed of 162 girls, distances of these magnitudes are to be expected by chance alone. Nonetheless, it is useful to identify the 2 girls whose  $d_i$  yielded  $p$ -values less than 0.01. They correspond to girls with subject ID = 128 and subject ID = 79. Recall that the former was identified as having extreme observations prior to menarche, with percent body fat increasing from 27% to 44% in a one-year interval. The other girl, with subject ID = 79, is unusual because she displayed a sudden decrease in percent body fat, from approximately 37% two years prior to menarche to 20%

around the time of menarche, and then maintained that level of percent body fat during the 3 years post-menarche. This pattern of a drop in percent body fat prior to menarche, coupled with almost no gain post-menarche, is at odds with the general pattern of change in the population.

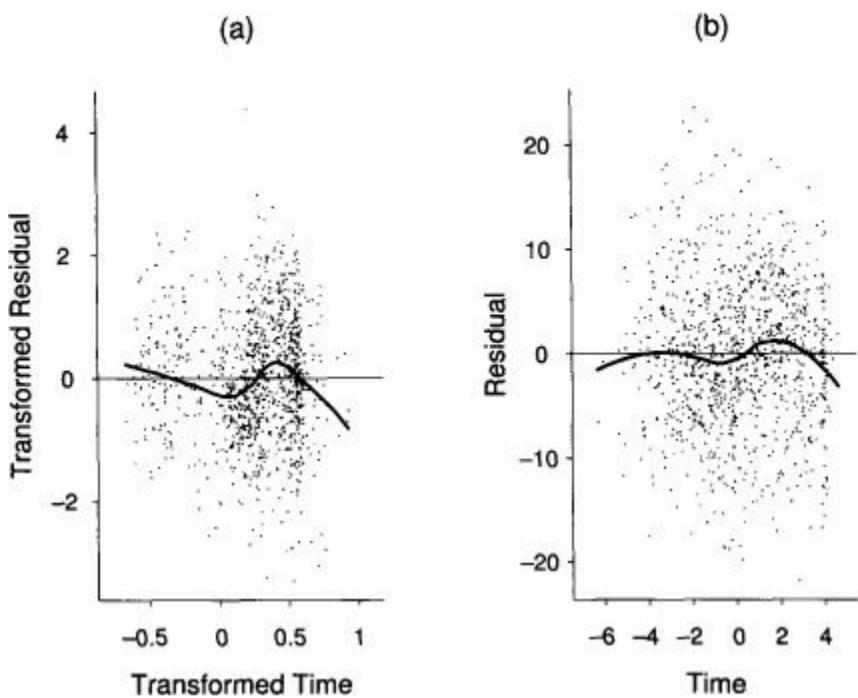
Next we consider scatterplots of the transformed and untransformed residuals versus the transformed and untransformed predicted values respectively. The scatterplots of the residuals in [Figure 10.3](#) display no obvious systematic pattern, with a random scatter around a constant mean of zero. However, when lowess smoothed curves are superimposed on the scatterplots, they do reveal some apparent curvature. Focusing on the transformed residuals, there appears to be a quadratic trend, although the fall in the lowess curve at the largest values of the transformed predicted values should be cautiously interpreted as the fitted curve is based on few observations at the extremities and is therefore likely to be unreliable in that region.

**Fig. 10.3** Scatterplot of (a) the transformed residuals versus transformed predicted values, and (b) the untransformed residuals versus predicted values, for the percent body fat data.



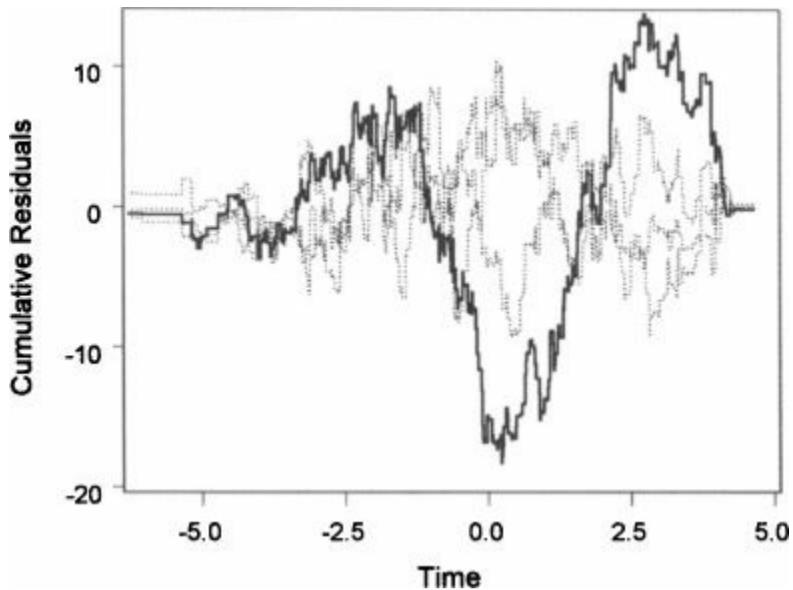
Because of the suggestion of curvature in [Figure 10.3](#), we next examine scatterplots of the (transformed) residuals versus (transformed) time (see [Figure 10.4](#)). These scatterplots of the transformed and untransformed residuals suggest curvature at (untransformed) times corresponding to approximately 2 to 4 years post-menarche. The pattern is more apparent in the scatterplot of the transformed residuals and can no longer be discounted due to sparseness of the observations at the extremities; this is the first pair of plots in which the transformed and untransformed data give a different impression. The curvature in the scatterplots suggests that the model for the mean response might be improved by the inclusion of a quadratic trend in the post-menarcheal period.

**Fig. 10.4** Scatterplot of (a) the transformed residuals versus transformed time, and (b) the untransformed residuals versus time, for the percent body fat data.

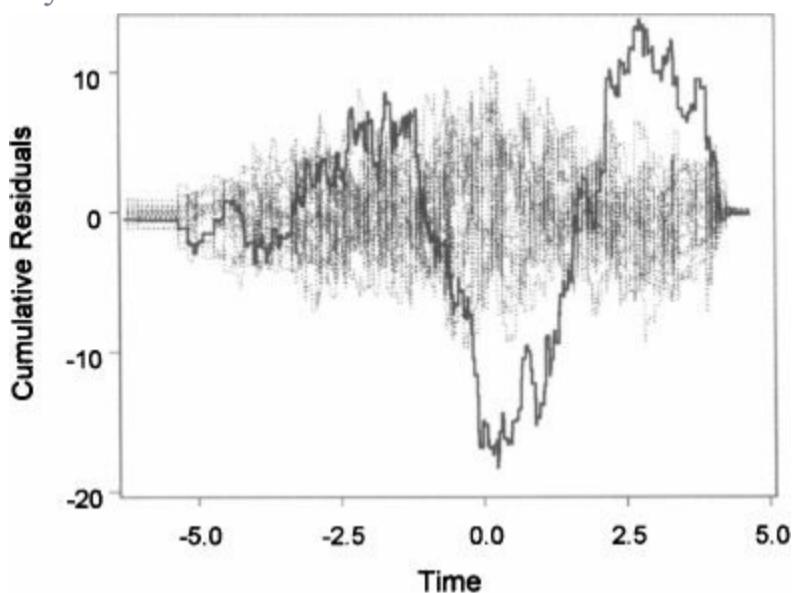


To remove some degree of subjectivity from our assessment of these graphical displays of the residuals, we consider the use of cumulative sums of residuals to assess the adequacy of the model. [Figure 10.5](#) shows a plot of the observed cumulative sum of the residuals (solid curve), with respect to time (relative to age of menarche). On the vertical axis is the cumulative sum of residuals; the horizontal axis denotes time (in years). Superimposed on the graph are three realizations (dotted curves) from the null distribution under the assumption that the model for the mean response is correctly specified. The three realizations of the cumulative sum under the null are computer simulated from the appropriate Gaussian mean-zero process. By comparing the observed cumulative sum to many different realizations under the null, it is possible to determine whether any apparent trend is systematic or due to chance fluctuations. From [Figure 10.5](#), the simulated realizations produce curves that appear to be closer to and intersect zero more often than the observed curve. By generating many more such realizations from the null distribution, it is possible to get both a graphical and numerical indication of whether the curve describing the observed cumulative sum displays a systematic pattern or simply natural variation. [Figure 10.6](#) shows a plot of the observed cumulative sum of the residuals and 20 realizations from the null distribution. It would appear that the observed cumulative sum displays a systematic pattern. In particular, the observed cumulative sum is too small in the 12 months after menarche (years 0 to 1) and too large 2 to 4 years after menarche. This suggests that the assumed functional form for time, in particular, after menarche, may not be adequate. This graphical assessment of fit can be complemented by a numerical assessment. Specifically, we obtained a  $p$ -value for the supremum test based on 10,000 simulated realizations from the null distribution. The maximum absolute value of the observed cumulative sum is 18.28. The supremum test yields a  $p$ -value of 0.0002, based on the 10,000 simulated realizations of the process under the null. That is, out of 10,000 simulated realizations, only two had a maximum absolute value that exceeded 18.28. Thus both the graphical and numerical results suggest that the functional form for time, in particular, after menarche (time = 0), may not be appropriate.

**Fig. 10.5** Plot of observed cumulative sum of residuals versus time since menarche (solid curve) and 3 simulated realizations (dotted curves) from the null distribution assuming a correctly specified model for mean percent body fat.

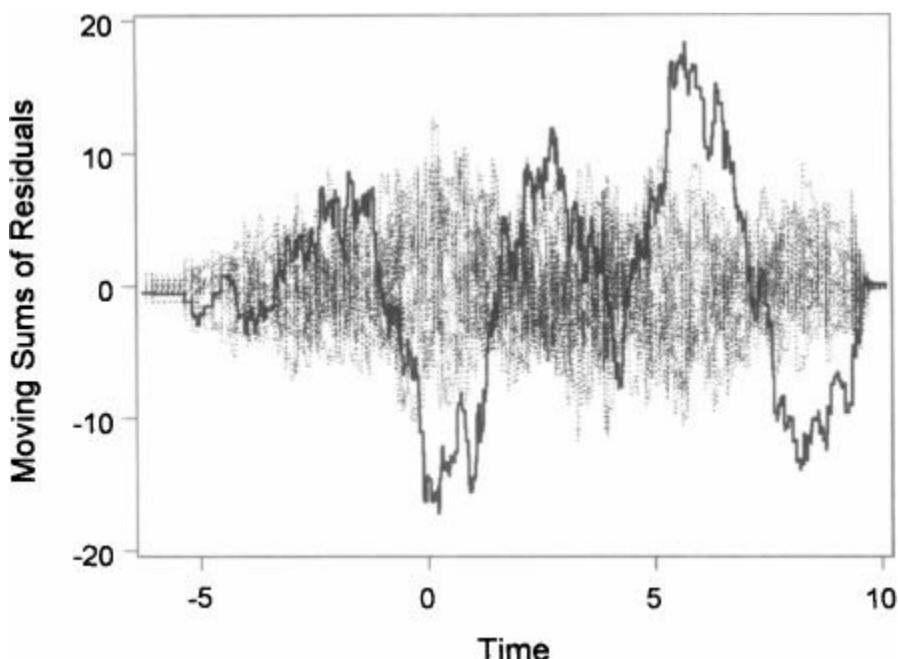


**Fig. 10.6** Plot of observed cumulative sum of residuals versus time since menarche (solid curve) and 20 simulated realizations (dotted curves) from the null distribution assuming a correctly specified model for mean percent body fat.



A similar plot can be constructed based on a moving sum rather than a cumulative sum. [Figure 10.7](#) shows a plot of the observed moving sum of the residuals, with block size equal to half the range of time (approximately 5.5 years). The observed curve in [Figure 10.7](#) also suggests that the moving sum of the residuals is too small in years 0-1 and too large in later years. In a similar fashion, we can complement this graphical display with a numerical assessment. The supremum test yields a  $p$ -value equal to 0.0004 (based on 10,000 simulated realizations), suggesting that the functional form for time may be inappropriate.

**Fig. 10.7** Plot of observed moving sum of residuals versus time since menarche (solid curve) and 20 simulated realizations (dotted curves) from the null distribution assuming a correctly specified model for mean percent body fat.



Next we consider a refinement to the model to allow for a quadratic trend in the post-menarcheal period. In particular, we assume that each girl has a piecewise linear-quadratic growth curve with a knot at the time of menarche and fit the following linear mixed effects model:

$$E(Y_{ij}|b_i) = \beta_1 + \beta_2 t_{ij} + \beta_3(t_{ij})_+ + \beta_4(t_{ij})_+^2 + b_{1i} + b_{2i}t_{ij} + b_{3i}(t_{ij})_+ + b_{4i}(t_{ij})_+^2,$$

where  $(t_{ij})_+^2 = t_{ij}^2$  if  $t_{ij} > 0$  and  $(t_{ij})_+^2 = 0$  if  $t_{ij} \leq 0$ . In this model each girl has a separate growth curve that can be described in terms of a linear trend for changes in response before menarche, and a quadratic trend for changes in response after menarche.

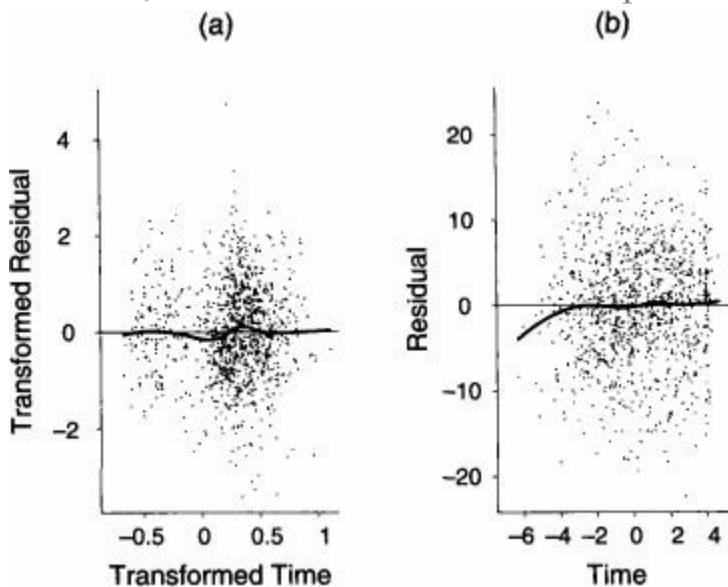
The REML estimates of the fixed effects,  $\beta = (\beta_1, \beta_2, \beta_3, \beta_4)'$ , are displayed in [Table 10.2](#). These results suggest that there is significant non-linearity in the post-menarcheal trend. The estimate of  $\beta_4$  indicates that increases in percent body fat are greatest around the time of menarche but level off at approximately 4 years following the onset of menarche. The results also suggest that there is no significant increase in percent body fat during the 3 to 4 years prior to menarche.

**Table 10.2** Estimated regression coefficients (fixed effects) and standard errors for the piecewise linear-quadratic model for the percent body fat data.

Variable	Estimate	SE	Z
Intercept	20.4201	0.5817	35.10
Time	-0.0155	0.1612	-0.10
(Time) <sub>+</sub>	4.8439	0.4055	11.94
(Time) <sub>+</sub> <sup>2</sup>	-0.6469	0.0772	-8.38

For this revised model we consider scatterplots of the (transformed) residuals versus (transformed) time (see [Figure 10.8](#)). The scatterplots of the transformed and untransformed residuals do not reveal any obvious systematic trends. When lowess smoothed curves are superimposed on the scatterplots, the curvature that was apparent in [Figure 10.4\(a\)](#) is no longer discernible in [Figure 10.8\(a\)](#). The apparent curvature in [Figure 10.8\(b\)](#), at approximately 5 to 6 years prior to menarche, can be discounted because the fitted curve is based on so few observations at this extremity and is therefore likely to be unreliable in that region. The inclusion of a quadratic trend in the post-menarcheal period has led to an improvement in fit as determined by both the Wald test for the quadratic trend ( $Z = -8.38, p < 0.0001$ ) and the examination of residual diagnostics.

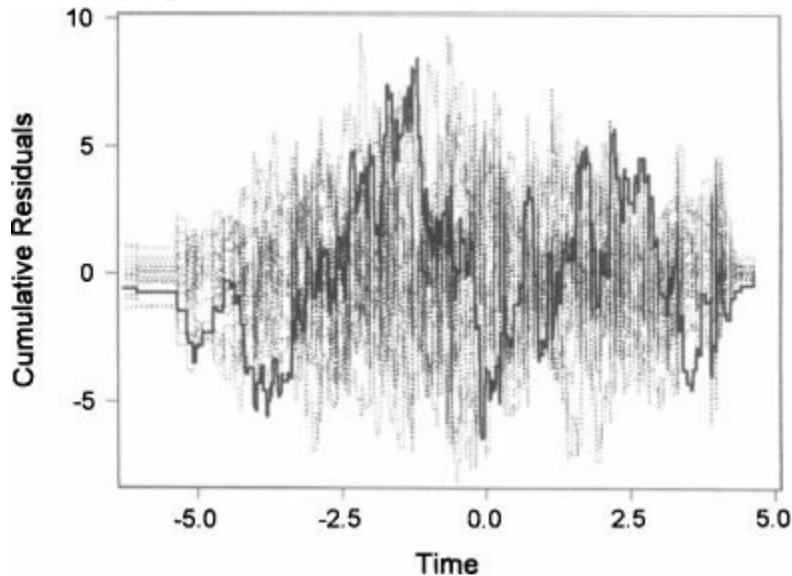
**Fig. 10.8** Scatterplot of (a) the transformed residuals versus transformed time, and (b) the untransformed residuals versus time, for the revised model for the percent body fat data.



Similarly we can assess the adequacy of the quadratic trend model using cumulative and moving sums of residuals. [Figure 10.9](#) shows a plot of the observed cumulative sum of the residuals, with respect to the covariate time; superimposed on the graph are 20 realizations from the Gaussian mean-zero null distribution. This plot suggests there is no systematic trend in the observed curve. This is confirmed by a numerical assessment. The maximum absolute value of the observed cumulative sum is 8.46, with corresponding  $p$ -value for the supremum test equal to 0.174. A similar plot can be constructed based on a moving sum and yields the same conclusion. Thus both the graphical and

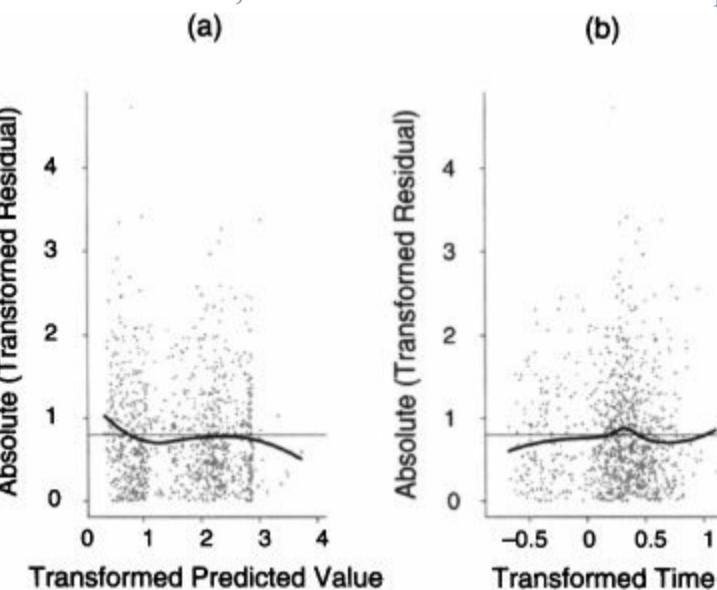
numerical results from the aggregation of the residuals suggest that the functional form for time (i.e., piecewise linear-quadratic) is adequate for these data.

**Fig. 10.9** Plot of observed cumulative sum of residuals versus time since menarche (solid curve) and 20 simulated realizations (dotted curves) from the null distribution assuming the revised model for mean percent body fat is correctly specified.



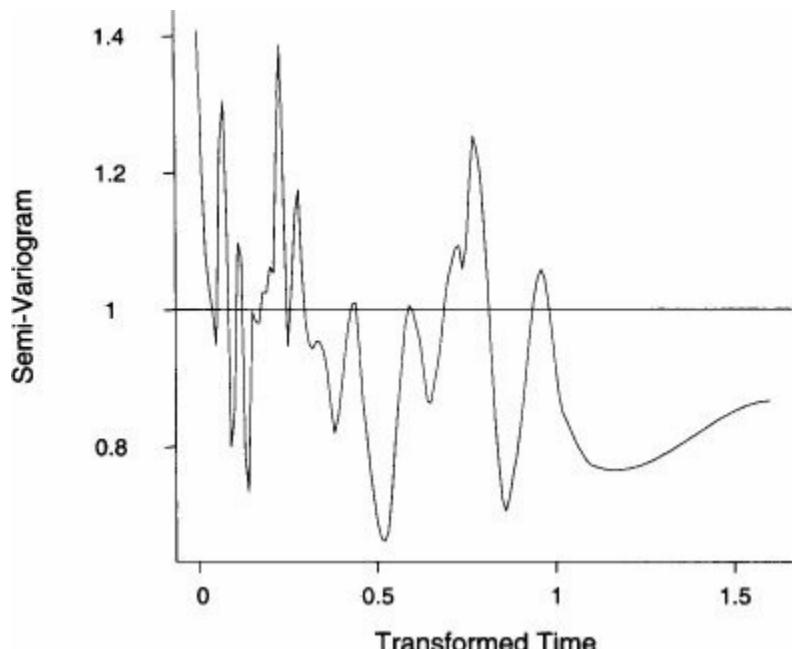
So far we have considered residual diagnostics for assessing the goodness of fit of the assumed model for the mean response. We can also assess the adequacy of the model for the variance by constructing a scatterplot of the absolute values of the transformed residuals,  $|\hat{r}_{ij}|$ , versus the transformed predicted values and transformed time (see [Figure 10.10](#)). The scatterplots in [Figure 10.10](#) indicate that there is no obvious systematic trend. When lowess smoothed curves are superimposed on the scatterplots, there is no evidence of a discernible departure from a straight line centered at approximately 0.8. Recall that if the transformed residuals are normal, with mean zero and unit variance, then the mean of the absolute values of the transformed residuals is 0.798. Thus we conclude that the variability is approximately constant for varying (transformed) predicted values and times.

**Fig. 10.10** Scatterplot of the absolute value of the transformed residuals versus (a) transformed predicted values, and (b) transformed time, for the revised model for the percent body fat data.



Finally, the adequacy of the overall model for the covariance matrix can be assessed by examining the empirical semi-variogram for the transformed residuals (see [Figure 10.11](#)). The empirical semi-variogram, estimated by fitting a lowess smoothed curve, appears to fluctuate randomly around the horizontal line centered at 1; it does not display any obvious systematic trend over time. This suggests that the assumed random effects structure for the covariance matrix is adequate for these data. For illustrative purposes we consider the empirical semi-variogram for the transformed residuals from the following mixed effects model:

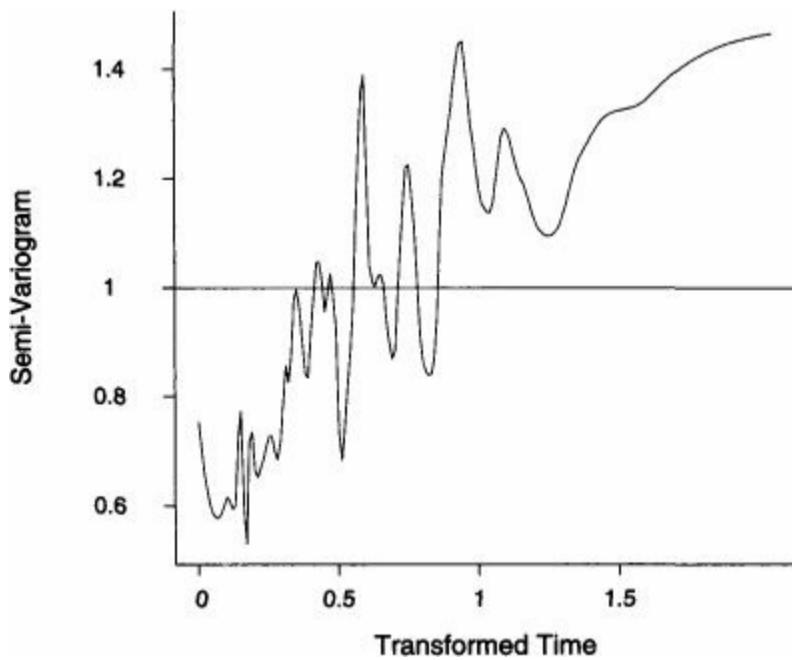
**Fig. 10.11** Empirical semi-variogram, estimated by fitting a lowess smoothed curve, for transformed residuals obtained from the revised model for the percent body fat data.



$$E(Y_{ij}|b_i) = \beta_1 + \beta_2 t_{ij} + \beta_3 (t_{ij})_+ + \beta_4 (t_{ij})_+^2 + b_{1i}.$$

This model includes only a single random effect,  $b_{1i}$ , allowing for heterogeneity in the intercepts (corresponding to percent body fat at menarche) but assuming no between-subject variation in the trajectories before or after menarche. The model retains the same set of fixed effects as before (i.e., the same model for the marginal means) but makes the strong assumption that the covariance matrix has a compound symmetry structure (i.e., a random intercepts only model). The adequacy of the compound symmetry model for the covariance matrix can be assessed by examining the empirical semi-variogram for the transformed residuals estimated from this model. [Figure 10.12](#) displays the empirical semi-variogram, estimated by fitting a lowess smoothed curve to the transformed residuals. The empirical semi-variogram is no longer centered around 1 and displays an increasing trend with time. This indicates that the compound symmetry assumption of constant variance and constant correlation is not tenable for these data.

**Fig. 10.12** Empirical semi-variogram, estimated by fitting a lowess smoothed curve, for transformed residuals obtained from the revised model for the fixed effects and with a compound symmetry assumption for the covariance matrix.



## 10.7 SUMMARY

In this chapter we have seen that many of the standard techniques for residual diagnostics and outlier detection in the univariate regression setting can be readily extended to longitudinal data. Simple residuals, based on the **observed minus predicted responses** at each occasion, are straightforward to calculate and can be produced by most statistical software packages for analyzing longitudinal data. Although these residuals have some shortcomings, they are probably adequate for most practical purposes. Systematic departures from model assumptions should be revealed by standard residual diagnostic plots of these simple residuals.

In Section 10.3 we discussed a particular transformation of the residuals for longitudinal data that makes them mimic residuals from a standard linear regression. The transformation is based on the Cholesky decomposition of the covariance matrix. Given an estimate of the covariance among the residuals, the Cholesky decomposition of the covariance matrix can be implemented in many standard statistical software packages. The transformed residuals can also be produced as standard output from some statistical packages (e.g., using the “normalized” residuals option with the `lme` function in the `n1me` package in R and S-Plus and the `VCIRY` option with PROC MIXED in SAS).

In Section 10.4 we discussed aggregating residuals by forming either “cumulative sums” or “moving sums” of residuals. The chief advantage of aggregating residuals is that a reference distribution is available that allows us to determine whether any apparent pattern in the *observed* sums of the residuals is evidence of a systematic trend or simply due to natural variation. This removes a large degree of subjectivity from the assessment of graphical displays of residuals. The aggregation of residuals can also be applied to the models for discrete longitudinal data discussed in Chapter 13 (see Section 13.3). The method has recently been implemented in PROC GENMOD in SAS.

Residual diagnostics for assessing the adequacy of the model for the covariance among repeated measurements are less well developed. The adequacy of the variance assumption can be informally assessed by examining a scatterplot of the absolute values of the transformed residuals versus the predicted values and/or time. A check on the adequacy of the overall model for the covariance is provided by a smoothed plot of the empirical semi-variogram.

Finally, residual analysis is useful for detecting outlying observations. The detection of outlying observations is important because these outliers can potentially have an inordinate influence on the analysis. Outliers require further investigation to ensure that they are not due to recording errors. When it has been established that the outliers are not the result of recording errors or other types of mistakes, then it can be useful to replicate the analysis with and without the outlying observations to determine their influence on substantive conclusions. However, we do not recommend the automatic exclusion or down-weighting of outliers. Residuals can also be used to identify outlying individuals who have unusual response profiles.

## **10.8 FURTHER READING**

Methods for residual analysis in standard linear regression for independent observations are well developed; a comprehensive description of techniques for residual analysis can be found in Cook and Weisberg (1982), and in many standard textbooks on linear regression. A discussion of residual diagnostics and the use of the semi-variogram for longitudinal data can be found in the review article by Laird et al. (1992).

# Bibliographic Notes

A discussion of the generalization of residual diagnostics to longitudinal data can be found in articles by Waternaux *et al.* (1989) and Waternaux and Ware (1991).

Lin, Wei, and Ying (2002) proposed the method of taking cumulative sums of residuals (and other aggregates of residuals) over covariates or predicted values to assess the adequacy of regression models for the mean response. They showed that the null distribution of the cumulative sum of residuals behaves approximately like a zero-mean Gaussian process whose realizations can be generated by computer simulation.

Historically the semi-variogram has been widely used in spatial statistics to represent the covariance structure in geostatistical data. The use of the semi-variogram for longitudinal data is described in detail in Chapter 10 (Section 10.4) of Verbeke and Molenberghs (2000) and Chapter 3 (Section 3.4) of Diggle *et al.* (2002); also see Chapters 2 (Section 2.5) and 5 (Section 5.4) of Diggle (1990).

## Problems

**10.1** In Section 8.8 we presented the results of analyses of a subset of the pulmonary function data collected in the Six Cities Study of Air Pollution and Health (Dockery *et al.*, 1983). The data consist of measurements of  $\text{FEV}_1$ , height, and age obtained from a randomly selected subset of the female participants living in Topeka, Kansas. Specifically, we considered the following model for  $\log(\text{FEV}_1)$ :

$$E(Y_{ij}|b_i) = \beta_1 + \beta_2 \text{Age}_{ij} + \beta_3 \log(\text{Ht})_{ij} + \beta_4 \text{Age}_{i1} + \beta_5 \log(\text{Ht})_{i1} + b_{1i} + b_{2i} \text{Age}_{ij},$$

where  $Y_{ij}$  is the  $\log(\text{FEV}_1)$  for the  $i^{th}$  child at the  $j^{th}$  visit,  $\text{Age}_{i1}$  and  $\log(\text{Ht})_{i1}$  are the initial or baseline age and  $\log(\text{height})$  for the  $i^{th}$  child. In this exercise we examine the residuals from the fitted model to assess the overall adequacy of the model for the mean response.

The raw data are stored in an external file: `fev1.dat`

Each row of the data set contains the following six variables:

ID Height Age Initial Height Initial Age  $\log(\text{FEV}_1)$

**10.1.1** Calculate the untransformed residuals,

$$\tau_{ij} = Y_{ij} - \hat{\mu}_{ij},$$

from the fitted model given above, where

$$\hat{\mu}_{ij} = \hat{\beta}_1 + \hat{\beta}_2 \text{Age}_{ij} + \hat{\beta}_3 \log(\text{Ht})_{ij} + \hat{\beta}_4 \text{Age}_{i1} + \hat{\beta}_5 \log(\text{Ht})_{i1}.$$

**10.1.2** Construct a histogram of the residuals. Comment on the shape of the distribution of the residuals.

**10.1.3** Construct a normal quantile plot (Q-Q plot) of the residuals. Does the plot display any systematic departures from a straight line? Does the plot suggest any potential outlying observations?

**10.1.4** On a single graph, construct a scatterplot of the residuals versus the predicted values and superimpose a lowess smoothed curve on the scatterplot. Does the plot display any systematic pattern?

**10.1.5** On a single graph, construct a scatterplot of the residuals versus age and superimpose a lowess smoothed curve on the scatterplot. Does the plot display any systematic pattern?

**10.1.6** On a single graph, construct a scatterplot of the residuals versus  $\log(\text{height})$  and superimpose a lowess smoothed curve on the scatterplot. Does the plot display any systematic pattern?

**10.1.7** On the basis of the residual diagnostics from Problems 10.1.2 through 10.1.6, comment on the overall adequacy of the model for the mean response. Can you suggest how the model for the mean response might be improved?

## *Part III*

# *Generalized Linear Models for Longitudinal Data*

# *Chapter 11*

## *Review of Generalized Linear Models*

### **11.1 INTRODUCTION**

In Part II we considered methods for analyzing longitudinal data when the response variable is continuous. In many biomedical applications the longitudinal response is not continuous, for example, the presence or absence of respiratory illness, or counts of the number of epileptic seizures in a four-week interval. When the longitudinal response is discrete (e.g., binary, ordinal, or a count), the linear models discussed in Part II are no longer appropriate for relating changes in the mean response to covariates. Instead, we consider extensions of *generalized linear models* for analyzing discrete longitudinal data.

Generalized linear models provide a unified class of models for regression analysis of independent observations of a discrete or continuous response. A straightforward application of generalized linear models to longitudinal data is not appropriate due to the correlation (or lack of independence) among observations obtained from the same individual. Instead, we consider extensions of this broad class of models to handle longitudinal responses. There are many ways to extend generalized linear models to account for the correlation among longitudinal observations, we consider two general, but quite distinct, approaches in Chapters 12 through 16.

In Chapters 12 and 13 we present a unified methodology for analyzing longitudinal data when the response variable is discrete or continuous. It does not require distributional assumptions for the observations, only a regression model for the mean response. That is, we describe a general method for analyzing diverse types of longitudinal responses that avoids making assumptions about the distribution of the vector of responses; the method relies solely on assumptions about how the mean response is related to covariates. Recall that in previous chapters we noted that the multivariate normal assumption was not so crucial in longitudinal analysis of a continuous response, provided the number of subjects is relatively large in comparison to the number of repeated measures and any missing data are MCAR. In Chapter 13 we provide some rationale for why the distributional assumption for the vector of responses can be relaxed. In Chapters 14 and 15 we consider an alternative extension of generalized linear models that accounts for the correlation among longitudinal data via the introduction of random effects. These models extend in a natural way the conceptual approach represented by the linear mixed effects models discussed in Chapter 8. In Part III we focus primarily on longitudinal analysis of a discrete response, although the general methodology described in these chapters can be applied equally to continuous responses.

A characteristic feature of generalized linear models is that a suitable non-linear transformation of the mean response is related to a linear function of the covariates. This non-linearity raises some additional issues concerning the interpretation of the regression coefficients in models for longitudinal data. In Chapter 16 we emphasize that different approaches for accounting for the source of within-subject association in longitudinal data can lead to models having regression coefficients with quite distinct interpretations. As a result, for the same data, there will be differences between the estimated regression coefficients obtained from the two distinct classes of models described in Chapters 12 through 15. In general, the choice among different classes of models for discrete longitudinal data must be made on subject-matter grounds.

One of the underlying themes that will be emphasized in Part III is that different models for discrete longitudinal data have somewhat different targets of inference. Thus, to ensure that the regression model parameters bear directly on the question of scientific interest, somewhat greater care is needed in the choice of model for discrete longitudinal data.

## 11.2 SALIENT FEATURES OF GENERALIZED LINEAR MODELS

In this section we provide a non-technical summary of the most salient features of generalized linear models for a single, univariate response. In later chapters we discuss how generalized linear models can be extended to handle longitudinal responses. A good grasp of the material in this section is all that is required for an understanding of the methodology for longitudinal data that will be described in subsequent chapters. In Section 11.7 we present a detailed and somewhat more technical overview of generalized linear models. Many of our readers, in particular, those encountering this topic for the first time, may find the material in Section 11.7 challenging. While we encourage all of our readers to skim through Section 11.7, we note that it can be omitted without loss of continuity.

Generalized linear models provide a unified method for analyzing diverse types of univariate responses (e.g., continuous, binary, ordinal, and count data). Generalized linear models are actually a broad class or collection of regression models, and they include as special cases the standard linear regression and analysis of variance (ANOVA) models for a normally distributed continuous response, logistic regression models for a binary or dichotomous response, and log-linear or Poisson regression models for counts. Although generalized linear models encompass a much broader range of regression models, these three are among the most widely used regression models in biomedical research. In this section we focus primarily on generalized linear models for binary and count data since, with the exception of continuous responses, these two data types are by far the most commonly encountered in applications. In Section 11.4 we discuss generalized linear models for ordinal data; we devote a separate section to ordinal regression models because, when regarded as generalized linear models, the models have certain non-standard features.

# Notation

Throughout this chapter we assume that we have  $N$  independent observations of a response variable,  $Y$ . We let  $Y_i$  ( $i = 1, \dots, N$ ) denote the response variable for the  $i^{th}$  subject. Associated with each response,  $Y_i$ , is a  $p \times 1$  vector of covariates,

$$X_i = \begin{pmatrix} X_{i1} \\ \vdots \\ X_{ip} \end{pmatrix}, \quad i = 1, \dots, N;$$

where  $X_{ik}$  denotes the  $k^{th}$  covariate for the  $i^{th}$  subject. Typically, although not always,  $X_{i1} = 1$  for all  $i$ , and then  $\beta_1$  is the intercept term in the regression model. Generalized linear models extend the standard linear regression model in a number of important ways, while also retaining some of its distinctive features. In particular, a generalized linear model for  $Y_i$  has the following three-part specification:

1. a distributional assumption,
2. a systematic component, and
3. a link function.

We consider each of these three components in turn.

# Distributional Assumption

Generalized linear models extend many of the basic concepts and ideas of standard linear regression analysis to settings where the response variable is discrete and can no longer be assumed to have a normal distribution. In particular, they extend the class of probability distributions for the response to include many of the distributions commonly used for modeling discrete responses. Generalized linear models assume that the response variable has a probability distribution belonging to the *exponential family* of distributions. The exponential family includes many distributions that the reader may already have encountered. For example, the normal, Bernoulli, binomial, and Poisson distributions all belong to the exponential family. The first component of a generalized linear model, the distributional assumption, specifies the *random component* of the model. That is, it specifies a probabilistic mechanism by which the responses are assumed to be generated.

Because the normal, Bernoulli, binomial, and Poisson distributions are members of the same family, they share some common statistical properties. In particular, the variance of the response can be expressed in terms of the product of a single scale or dispersion parameter,  $\phi$ , and a *variance function*, denoted  $v(\mu_i)$ ; the latter being a known function of the mean,  $\mu_i$ . That is,

$$\text{Var}(Y_i) = \phi v(\mu_i),$$

where  $\phi > 0$ . The variance function,  $v(\mu_i)$ , describes how the variance is functionally related to the mean of the response. The variance functions for the normal, Bernoulli, and Poisson distributions are summarized in [Table 11.1](#). For many distributions for discrete data,  $\phi$  is not a parameter that requires estimation but is a known constant (e.g.,  $\phi = 1$  for the Bernoulli and Poisson distributions); for other distributions  $\phi$  is an unknown parameter (e.g.,  $\phi$  is the variance of the normal distribution).

**Table 11.1** Canonical link and variance functions for the normal, Bernoulli, and Poisson distributions.

Distribution	Variance Function, $v(\mu)$	Canonical Link
Normal	$v(\mu) = 1$	Identity: $\mu = \eta$
Bernoulli	$v(\mu) = \mu(1 - \mu)$	Logit: $\log\left(\frac{\mu}{1-\mu}\right) = \eta$
Poisson	$v(\mu) = \mu$	Log: $\log(\mu) = \eta$

For the Bernoulli and Poisson distributions, the variance depends on the mean. This dependence of the variance on the mean is a characteristic feature of most distributions for discrete responses. On the other hand, for the normal distribution, the variance does not depend on the mean; that is,  $\text{Var}(Y_i) = \phi$  (and the variance function,  $v(\mu_i) = 1$ ). This provides some rationale for why the assumption of homogeneity of variance (or common variance) is generally adopted in the standard linear regression model for normally distributed responses. In some applications, however, the homogeneity of variance assumption is too restrictive and the variance may depend on covariates. In later sections we briefly mention how restrictive assumptions about the variance of  $Y_i$  can be relaxed.

# Systematic Component

Generalized linear models not only share a common family of distributions, they also share a common regression formulation. An important aspect of the standard linear regression model that is retained in all generalized linear models is the linear regression component. This is the *systematic component* of a generalized linear model, and it specifies that the effects of the covariates,  $X_i$ , on the mean of  $Y_i$  can be expressed in terms of the following *linear predictor*, denoted by  $\eta_i$ ,

$$\eta_i = \beta_1 X_{i1} + \beta_2 X_{i2} + \cdots + \beta_p X_{ip},$$

where, typically,  $X_{i1} = 1$  for all  $i$ , and then  $\beta_1$  is the intercept. The linear predictor is simply a linear combination of the unknown regression coefficients,  $\beta = (\beta_1, \dots, \beta_p)'$  and the covariates,  $X_i$ .

The key word here is *linear*. The term “linear” in generalized linear models means that  $\eta_i$  must be linear in the regression parameters. This implies that the mean response (or any transformation of the mean response) can be expressed as a simple weighted sum of the regression parameters,  $\beta$ . For example,

$$\eta_i = \beta_1 + \beta_2 X_i,$$

$$\eta_i = \beta_1 + \beta_2 \log(X_i),$$

and

$$\eta_i = \beta_1 + \beta_2 X_i + \beta_3 X_i^2,$$

are all cases where  $\eta_i$  is linear in the regression coefficients, even if it is non-linear in  $X_i$ . However,

$$\eta_i = \beta_1 + e^{\beta_2 X_i},$$

and

$$\eta_i = \beta_1 / (1 + \beta_2 e^{-\beta_3 X_i})$$

are examples where  $\eta_i$  is not linear in the regression parameters and the latter types of non-linearities are not included in the class of generalized linear models.

Thus the linearity strictly applies to the regression parameters,  $\beta$ , but not necessarily to the covariates. As a result the linearity restriction does not preclude relationships between the mean response and the covariates that are non-linear. This latter type of non-linearity is easily accommodated by taking appropriate transformations of the mean response (see below) and/or by transformation of the covariates (e.g.,  $\log(X)$  or  $X^2$ ).

# Link Function

The final way in which generalized linear models extend the standard linear regression model is by taking a suitable transformation of the mean response and relating the transformed mean response to the covariates. This is achieved by the introduction of a *link function*. The link function applies a transformation to the mean and then links the covariates, via the linear predictor, to the transformed mean of the distribution of the responses,

$$g(\mu_i) = \eta_i = \beta_1 X_{i1} + \beta_2 X_{i2} + \cdots + \beta_p X_{ip} = \sum_{k=1}^p \beta_k X_{ik} = X_i' \beta,$$

where the link function  $g(\cdot)$  is some known function, for example,  $\log(\mu_i)$ . This implies that it is the transformed mean response that changes linearly with changes in the values of the covariates.

Thus, while in the standard linear regression model the mean response is related directly to a linear combination of the covariates, in generalized linear models, it is some appropriate transformation of the mean response, for example,  $\log(\mu_i)$ , that is related to a linear combination of the covariates. The linearity applies to a transformation of the mean response, or, put in a somewhat different way, the effects of covariates are assumed to be additive on a suitably transformed scale for the mean response.

The use of non-linear link functions, for instance,  $\log(\mu_i)$ , ensures that the model produces predictions of the mean response that are within the allowable range. For example, when analyzing a binary response,  $\mu_i$  has interpretation in terms of the probability of “success” (with  $0 < \mu_i < 1$ ). If the mean response, here the probability of success, is related directly to a linear combination of the covariates, the model can yield predicted probabilities outside of the range from 0 to 1. The use of certain non-linear link functions ensures that this cannot happen, while at the same time allowing an unbounded range of values for the regression parameters,  $\beta$ .

We can distinguish two main types of link functions, *canonical* link functions and *non-canonical* link functions. The former are unique and can be derived for any selected distribution; the latter are somewhat arbitrary and bear no direct relation to the selected distribution. For example, the logit link function is the canonical link function associated with the Bernoulli and binomial distributions; the probit link function is a non-canonical link function for these distributions that is often adopted for the analysis of binary data from toxicological experiments. Although, in principle, any suitable link function can be used to relate the mean response to the covariates, the choice of a canonical link function produces many of the most widely used regression models. The canonical link functions for the normal, Bernoulli, and Poisson distributions are summarized in [Table 11.1](#).

In summary, in generalized linear models, the distribution of the response is assumed to belong to a single family of distributions known as the exponential family. The exponential family includes the normal, Bernoulli, binomial, and Poisson distributions. A transformation of the mean response is then linearly related to the covariates, via an appropriate link function. Because generalized linear models make distributional assumptions about the response variable, the regression parameters can be estimated using the method of maximum likelihood. The maximum likelihood estimates of the regression coefficients,  $\beta$ , are simply those values of  $\beta$  that are most probable (or most “likely”) for the data that have actually been observed. The method of maximum likelihood provides a very general technique for estimation and for inference, that is, for estimating  $\beta$ , constructing confidence intervals, testing hypotheses, and assessing the adequacy of models. All of these ideas will be elaborated in Section 11.3, where we focus on two special cases of generalized linear models: logistic regression for a binary response and log-linear regression for counts. Although generalized linear models provide a very broad and flexible collection of regression models for analyzing diverse types of responses, they do have one very important restriction: they assume that observations on the response variable are independent of one another. In later chapters we will discuss how this restriction to independent observations can be relaxed to accommodate the correlated nature of the responses arising from longitudinal studies.

## 11.3 ILLUSTRATIVE EXAMPLES

To clarify the main ideas presented in the previous sections, we consider in greater detail two special cases of generalized linear models: logistic regression for a binary response and log-linear regression for counts. We consider each of these models in terms of their three-part specification as a generalized linear model. We also emphasize the interpretation of the regression coefficients,  $\beta$ , in these models. In Section 11.4 we discuss generalized linear models for ordinal data. The description of methods for extending generalized linear models for longitudinal responses presented in later chapters will assume a good working knowledge of these important regression models. As a result the reader is encouraged to master the material in this section (and Section 11.4) before proceeding to Chapters 12 through 16. This section can be skimmed through for those with a strong background in logistic and log-linear regression models.

## 11.3.1 Logistic Regression for Binary Responses

Logistic regression is used widely to describe the relationship between a binary response variable (e.g., denoting “success” or “failure”) and a set of covariates. In common with standard linear regression, the primary objective of logistic regression is to relate the mean of the response to a set of covariates. However, the response variable is binary rather than continuous and this has a number of consequences for modeling the mean. In this section we describe the main features of logistic regression and highlight its three-part specification as a generalized linear model. We also consider various aspects of interpretation of logistic regression coefficients. An example, using data of low-birth-weight infants, is used to illustrate the main ideas.

Let  $Y_i$  denote a binary response variable, whose two categories, for convenience, are often referred to as “success” or “failure.” For example,  $Y_i$  might indicate the presence or absence of a disease. Denoting the two possible outcomes for  $Y_i$  by 1 (for “success”) and 0 (for “failure”), the probability distribution of  $Y_i$  is Bernoulli, with  $\Pr(Y_i = 1) = \mu_i$  (and correspondingly,  $\Pr(Y_i = 0) = 1 - \mu_i$ ). The primary goal of logistic regression is to describe the effects of changes in a set of covariates,  $X_i$ , on the mean  $\mu_i$ . For ease of exposition we first consider the simple case where there is only a single covariate,  $X_i$ . Generalizations to more than one covariate will be considered later.

Since the analytic goal is to investigate the relationship between  $\mu_i$  and  $X_i$ , and since linear regression plays such a dominant role in applications, it may at first seem natural to assume a linear model relating the mean of  $Y_i$  to  $X_i$ ,

$$E(Y_i|X_i) = \mu_i = \beta_1 + \beta_2 X_i.$$

However, this linear model for the probabilities has one obvious difficulty. Expressing  $\mu_i$  as a linear function violates the restriction that probabilities must lie within the range from 0 to 1. For sufficiently large or small values of  $X_i$ , this regression model will yield predicted probabilities outside of the range from 0 to 1. A further difficulty with the linear model for  $\mu_i$  is that we often expect a non-linear relationship between  $\mu_i$  and  $X_i$ . For example, a 0.2 unit increase in  $\mu_i$  might be considered more “extreme” when  $\mu_i = 0.1$  than when  $\mu_i = 0.5$ . In terms of ratios, the change from  $\mu_i = 0.1$  to  $\mu_i = 0.3$  represents a three-fold or 200% increase, whereas the change from  $\mu_i = 0.5$  to  $\mu_i = 0.7$  represents only a 40% increase. In a sense, the units of measurement for a probability (or proportion) are often not considered to be constant over the range from 0 to 1. The linear probability model simply does not take this into consideration when relating  $\mu_i$  to  $X_i$ . Note also that the usual assumption of homogeneity of variance (or constant variance) in linear regression would be violated since the variance of a binary response explicitly depends on the mean, with

$$\text{Var}(Y_i) = \mu_i(1 - \mu_i).$$

To circumvent these difficulties with the linear probability model, a non-linear transformation can be applied to  $\mu_i$  and the transformed probabilities are related linearly to  $X_i$ . When the logit or logistic function,  $\log\{\mu_i/(1 - \mu_i)\}$ , is adopted, the resulting model

$$\log\left(\frac{\mu_i}{1 - \mu_i}\right) = \text{logit}(\mu_i) = \beta_1 + \beta_2 X_i,$$

is known as the *logistic regression* model. Recall that if  $\mu_i$  is the probability of success, then  $\mu_i(1 - \mu_i)$  is known as the *odds* of success. For example, if the probability of success is 0.8 then the odds of success is 4 (or 0.8/0.2) to 1. That is, the probability of success is 4 times as large as the probability of failure. Thus the logistic regression model assumes a linear relationship between the log odds of success and  $X_i$ . For the reader unfamiliar with logistic regression, it is useful to bear in mind that the transformation of  $\mu_i$  in logistic regression has the following property: as the probability of success,  $\mu_i$ , increases, so too does the odds of success and the log odds of success; similarly, as the probability of success decreases, so too does the odds of success and the log odds of success.

Next consider the interpretation of the logistic regression coefficients,  $\beta_1$  and  $\beta_2$ . For the special

case where the predictor variable  $X_i$  is dichotomous, taking values of 0 and 1, the logistic regression slope,  $\beta_2$ , has a simple and very attractive interpretation in terms of the log odds ratio (comparing the log odds of success when  $X_i = 1$  to the log odds of success when  $X_i = 0$ ). That is,

$$\text{logit}(\mu_i|X_i = 1) - \text{logit}(\mu_i|X_i = 0) = (\beta_1 + \beta_2) - \beta_1 = \beta_2.$$

Thus  $\exp(\beta_2)$  has interpretation as the odds ratio of the response for the two possible values of the covariate.

In simple linear regression the interpretation of the slope of the regression is in terms of changes in the mean of  $Y_i$  for a single-unit change in  $X_i$ . Similarly, for arbitrary  $X_i$ , the logistic regression slope  $\beta_2$  has interpretation as the change in the log odds (of success) for a unit change in  $X_i$ . Equivalently, a unit change in  $X_i$  increases or decreases the odds of success *multiplicatively* by a factor of  $\exp(\beta_2)$ . Also recall that the intercept in simple linear regression has interpretation as the mean value of the response variable when  $X_i$  is equal to zero. Similarly, the logistic regression intercept,  $\beta_1$ , has interpretation as the log odds (of success) when  $X_i = 0$ ; alternatively,

$$\frac{\exp(\beta_1)}{1 + \exp(\beta_1)}$$

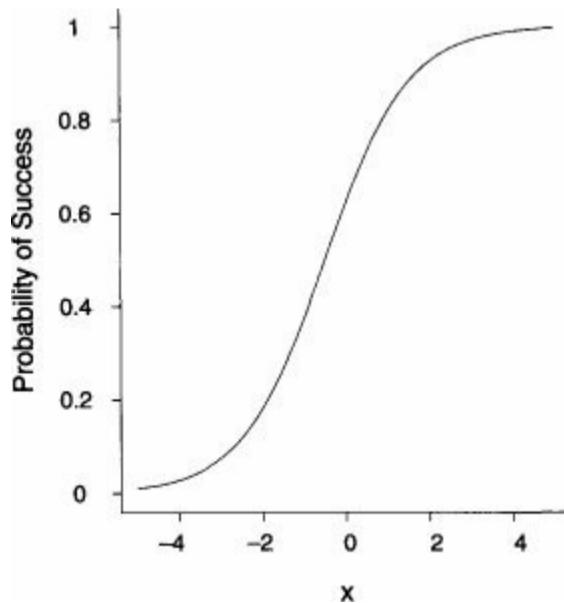
is the probability of success when  $X_i = 0$ .

The logistic regression model can also be expressed in terms of the probability of success,  $\mu_i$ ,

$$\mu_i = \frac{\exp(\beta_1 + \beta_2 X_i)}{1 + \exp(\beta_1 + \beta_2 X_i)}.$$

While the latter expression may appear to be somewhat more complicated, this is simply an equivalent way of expressing the logistic regression model. That is, logistic regression describes how the log odds,  $\log(\frac{\mu_i}{1-\mu_i})$ , has a linear relationship with  $X_i$ , which is equivalent to describing how  $\mu_i$  has a sigmoidal or S-shaped relationship with increasing values of  $\beta_1 + \beta_2 X_i$ . (See [Figure 11.1](#) for a plot of  $\mu$  versus  $X$  when  $\beta_1 = 0.5$  and  $\beta_2 = 0.9$ .) Of note, the logistic transformation ensures that the predicted probabilities are restricted to the range from 0 and 1, while allowing an unbounded range for  $\beta_1$  and  $\beta_2$ .

**Fig. 11.1** Plot of logistic response function, with success probability,  $\mu = \frac{e^{0.5+0.9X}}{1 + e^{0.5+0.9X}}$ .



When viewed as a generalized linear model, logistic regression is simply the special case where the distribution of  $Y_i$  is assumed to be Bernoulli (a member of the exponential family) and a logit link function, the canonical link function, has been adopted. Because the Bernoulli distribution is a one-parameter exponential family distribution, the variance of  $Y_i$  can be expressed explicitly in terms of the mean, via the following variance function:

$$\text{Var}(Y_i) = v(\mu_i) = \mu_i(1 - \mu_i)$$

and  $\phi = 1$ . For the Bernoulli distribution, the dispersion parameter is a fixed and known constant ( $\phi = 1$ ).

When  $X_i$  is a discrete covariate with  $J$  distinct categories or levels (e.g., treatment groups), the

binary responses for the  $N$  individuals can be grouped. Let  $m_j$  denote the number of individuals with the  $j^{th}$  covariate pattern, and let  $Y_j$  denote the number of successes among the  $m_j$  individuals, for  $j = 1, \dots, J$ . We may provisionally assume that all individuals within a group respond independently with constant probability of success,  $\mu_j$ , depending only on group. Then  $Y_j$ , the number of successes in the  $j^{th}$  group, has a binomial distribution with probability of success,  $\mu_j$ . The binomial distribution belongs to the exponential family and the probability of success,  $\mu_j$ , can be related to group using a logit link function. For the binomial distribution, the mean or expected number of successes for the  $j^{th}$  covariate pattern is

$$E(Y_j) = m_j \mu_j.$$

There is a well-known relationship between the mean and variance of  $Y_j$ , with

$$\text{Var}(Y_j) = m_j \mu_j (1 - \mu_j).$$

However, in many biomedical applications, counts of the number of successes have variability that far exceeds that predicted by the binomial distribution; this phenomenon is often referred to as *overdispersion* (although underdispersion can also arise, it is far less common). Overdispersion is a common indicator of failure of the binomial assumptions: independent observations with constant probability of success. That is, overdispersion can be represented either by a positive correlation between the responses or by variation in the response probabilities. To allow for overdispersion or extra-binomial variation, a scale factor  $\phi$  (with  $\phi \neq 1$ ) is often included in the specification of the binomial variance,

$$\text{Var}(Y_j) = \phi m_j \mu_j (1 - \mu_j).$$

Failure to account for overdispersion has negligible impact of the estimated logistic regression coefficients. That is, the regression parameter estimates are consistent and there is usually little loss of efficiency. Neglecting overdispersion, however, results in the standard errors being underestimated and failure to make an adjustment to the nominal standard errors can result in misleading inferences concerning the regression coefficients (e.g., confidence intervals that are too narrow and  $p$ -values that are too small). Also model selection strategies based on likelihood ratio tests or on information criteria, such as the Akaike information criterion (AIC), will perform poorly. When overdispersion is ignored, a model with too many parameters is likely to be selected and thus can lead to overinterpretation of these parameters (e.g., unnecessary inclusion of interactions). Adjustments to the nominal standard errors to account for overdispersion can be made either by including a scale factor  $\phi$  in the specification of the binomial variance,

$$\text{Var}(Y_j) = \phi m_j \mu_j (1 - \mu_j),$$

or by basing standard errors on the so-called “sandwich” estimator of  $\text{Cov}(\hat{\beta})$ ; the latter will be discussed in greater detail in Chapter 13.

So far we have only considered the simple case where there is a single predictor  $X_i$ . Next, we consider the case where  $X_i$  is a  $p \times 1$  vector of covariates. The logistic regression model becomes

$$\log\{\mu_i/(1 - \mu_i)\} = \beta_1 X_{i1} + \beta_2 X_{i2} + \dots + \beta_p X_{ip},$$

where  $X_{i1} = 1$  for all  $i = 1, \dots, N$ . The logistic regression coefficients in this model have the following interpretations. Each of the logistic regression slopes,  $\beta_k$  (for  $k = 2, \dots, p$ ), represents the change in the log odds (of success) for a unit change in  $X_{ik}$  given that all of the other predictor variables remain constant. This is completely analogous to the interpretation of the regression coefficients in multiple linear regression. Thus, by holding the remaining predictors at some fixed set of values and not allowing them to vary with any changes in  $X_{ik}$ , a single-unit increase in  $X_{ik}$  is predicted to increase or decrease the log odds of success by an amount  $\beta_k$ . Equivalently, a single-unit increase in  $X_{ik}$  increases or decreases the odds of success *multiplicatively* by a factor of  $\exp(\beta_k)$ . The logistic regression intercept,  $\beta_1$ , now has interpretation as the log odds (of success) when all covariate values are set to zero. Alternatively,

$$\frac{\exp(\beta_1)}{1 + \exp(\beta_1)}$$

is the probability of success when  $X_{i2} = X_{i3} = \dots = X_{ip} = 0$ .

Finally, the logistic regression model for binary data can also be developed through the notion of an underlying latent variable distribution. Suppose that  $L_i$  is a latent (i.e., unobserved) continuous variable and that a positive response is observed only when  $L_i$  exceeds some threshold denoted by  $\tau$ . The observed binary response can be thought of as a categorization of the unobservable latent variable, above and below the threshold  $\tau$ . That is,

$$Y_i = 1 \text{ if } L_i > \tau,$$

$$Y_i = 0 \text{ if } L_i \leq \tau.$$

Suppose that the latent variable,  $L_i$ , has a standard logistic distribution. The standard logistic distribution (with mean zero and variance  $\pi^2/3$ ) is a symmetric distribution and is very similar to the standard normal distribution, except that it has longer tails (and larger variance). Then, using calculus, it can be shown that the relationship between the observable binary response variable,  $Y_i$ , and the unobservable latent variable,  $L_i$ , is given by

$$\begin{aligned} \Pr(Y_i = 1) &= \Pr(L_i > \tau) \\ &= \int_{\tau}^{\infty} \frac{\exp(u)}{\{1 + \exp(u)\}^2} du \\ &= \frac{\exp(-\tau)}{1 + \exp(-\tau)}. \end{aligned}$$

Next suppose that the following linear model for  $L_i$  holds:

$$L_i = \beta_1 X_{i1} + \beta_2 X_{i2} + \cdots + \beta_p X_{ip} + e_i = X'_i \beta + e_i,$$

where  $e_i$  (or  $L_i - X'_i \beta$ ) is assumed to have a standard logistic distribution, with mean zero and variance  $\pi^2/3$ . Here we regard the threshold of  $L_i$  as fixed and the location or mean of the distribution of  $L_i$  as changing with  $X_i$ . Without loss of generality, we can assume the threshold for categorizing  $L_i$  is  $\tau = 0$ , since any non-zero values for the threshold would simply be absorbed into the intercept term in the linear predictor,  $X'_i \beta$ . Then the relationship between  $Y_i$  and  $L_i$  results in a logistic regression model for  $\Pr(Y_i = 1)$ . That is,

$$\begin{aligned} \Pr(Y_i = 1) &= \Pr(L_i > 0) \\ &= \Pr(L_i - X'_i \beta > -X'_i \beta) \\ &= \frac{\exp(X'_i \beta)}{1 + \exp(X'_i \beta)}. \end{aligned}$$

Thus the linear model for  $L_i$  with standard logistic errors,

$$L_i = X'_i \beta + e_i,$$

implies the logistic regression model for  $Y_i$ ,

$$\log \left\{ \frac{\Pr(Y_i = 1)}{1 - \Pr(Y_i = 1)} \right\} = X'_i \beta.$$

Similarly, if a probit link function is adopted instead of a logit link function, the linear model for  $L_i$  with standard normal errors, instead of standard logistic errors, implies a probit regression model for  $Y_i$ .

Although the logistic regression model can be derived from the notion of a latent variable distribution, assuming the existence of a latent variable is not a necessary requirement for the use of logistic regression models (indeed, in practice, the existence of a latent variable is usually not verifiable from the data). In later chapters we use the notion of an underlying latent variable distribution to derive analogues of the between-subject and within-subject sources of variability in models for longitudinal binary responses.

# Illustration

Next we consider an application of logistic regression to illustrate how the model can be used in practice. The data are from a study of low-birth-weight infants in a neonatal intensive care unit. In this example we are interested in the development of bronchopulmonary dysplasia (BPD), a chronic lung disease, in a sample of 223 infants weighing less than 1750 grams (Van Marter et al., 1990).

Let  $Y_i$  be a binary response, with  $Y_i = 1$  if the  $i^{th}$  infant develops BPD by day 28 of life and  $Y_i = 0$  otherwise (where BPD is defined by both oxygen requirement and compatible chest radiograph). To examine whether there is an association between the risk of BPD and birth weight (in grams  $\times 10^{-2}$ ), we consider the following logistic regression model:

$$\log \{\mu_i / (1 - \mu_i)\} = \beta_1 + \beta_2 \text{Weight}_i,$$

where  $\mu_i = E(Y_i) = \Pr(Y_i = 1)$ . For the 223 infants in the sample, the estimated logistic regression parameters (and standard errors), obtained using maximum likelihood, are displayed in [Table 11.2](#).

**Table 11.2** Estimated coefficients and standard errors for logistic regression of BPD on birth weight.

Variable	Estimate	SE	Z
Intercept	4.0343	0.6958	5.798
Birth Weight	-0.4229	0.0641	-6.599

The estimated logistic regression is

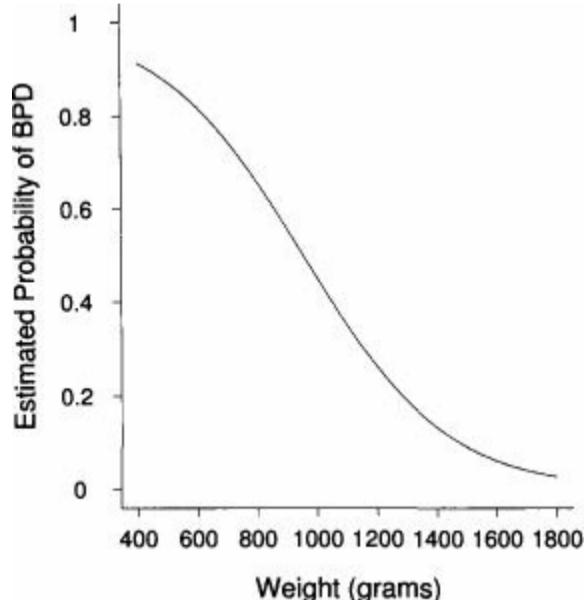
$$\log \{\hat{\mu}_i / (1 - \hat{\mu}_i)\} = 4.0343 - 0.4229 \text{Weight}_i.$$

When compared to its standard error, the ML estimate of  $\beta_2$ , the slope for birth weight, is significantly different from zero at the 0.05 level. The results from the logistic regression analysis indicate that the risk of BPD decreases with increasing birth weight. Specifically, the estimate of  $\beta_2$  implies that for every 100 gram increase in birth weight, the log odds of BPD decreases by 0.42. For example, the odds of BPD for an infant weighing 1200 grams (approximately 2.5 pounds) is

$$\exp(4.0343 - 12 \times 0.4229) = \exp(-1.0405) = 0.353.$$

Thus the predicted probability of BPD is  $0.353 / (1 + 0.353) \approx 0.26$ . The estimated probability of BPD can be calculated at any birth weight and a plot of the estimated probability versus weight produces the characteristic sigmoidal or S-shaped curve displayed in [Figure 11.2](#).

**Fig. 11.2** Plot of estimated logistic response function of BPD on birth weight based on a sample of 223 infants with birth weight less than 1750 grams.



Next suppose that we include two additional covariates, gestational age (in weeks) and presence of toxemia (with 1 denoting the presence of toxemia and 0 its absence). That is, we consider the following logistic regression model:

$$\log \{\mu_i / (1 - \mu_i)\} = \beta_1 + \beta_2 \text{Weight}_i + \beta_3 \text{Age}_i + \beta_4 \text{Toxemia}_i.$$

For the 223 infants in the sample, the estimated logistic regression parameters (and standard errors) are displayed in [Table 11.3](#).

**Table 11.3** Estimated coefficients and standard errors for logistic regression of BPD on birth weight, gestational age and toxemia.

Variable	Estimate	SE	Z
Intercept	13.9361	2.9826	4.672
Birth Weight	-0.2644	0.0812	-3.254
Gestational Age	-0.3885	0.1149	-3.382
Toxemia	-1.3437	0.6075	-2.212

The estimated coefficient for birth weight has now decreased, when gestational age and toxemia are included in the analysis. Nonetheless, the estimate of  $\beta_2$  remains significantly different from zero at the 0.05 level. The estimated coefficient for gestational age has interpretation in terms of the change in the log odds of BPD for a 1-week increase in gestational age, adjusting for birth weight and toxemia. Specifically, a 1-week increase in gestational age is associated with a 0.39 decrease in the log odds of developing BPD. Finally, the estimated coefficient for toxemia has interpretation in terms of the log odds ratio, comparing mothers who were diagnosed with toxemia to mothers who were not, while adjusting for the effects of birth weight and gestational age. Specifically, the adjusted odds ratio is 0.26 (or  $e^{-1.34}$ ) and indicates that infants of mothers diagnosed with toxemia have approximately a quarter the risk of developing BPD.

## 11.3.2 Log-Linear Regression for Counts

Log-linear regression, often referred to as Poisson regression, is used widely for the analysis of counts of the number of times some event occurs in either time or space. For example, the response variable might be the count of the number of epileptic seizures a particular patient experiences in a 4-week interval. Alternatively, the response might be a count of bacteria present in a fixed volume of bacterial suspension. In either case the response variable  $Y_i$  is a count, and the objective of log-linear regression is to relate the mean or expected count to a set of covariates.

If the occurrences of some event are counted within an interval of time (or sometimes volume or area), then the *rate* at which the event occurs is usually of more direct interest than the corresponding count. That is, the count or absolute number of events is often not satisfactory because any comparisons depend almost entirely on the “time at risk” (or, in other contexts, the sizes of the groups or areas of regions) that generated the observations. For example, it would not be very meaningful to compare counts of the number of seizures in a 4-week interval with counts of the number of seizures in a 12-month interval since it seems reasonable to suppose that the number of seizures is directly proportional to the period at risk. When the “time at risk” is not the same for all observations, a rate provides a meaningful basis for direct comparison. In either case the primary objective of log-linear regression is to relate the expected counts or rates to a set of covariates.

When the response is a count it is often reasonable to assume that  $Y_i$  has a Poisson distribution, although it is important to note that this is an assumption and it may not be valid. This is in contrast to the binary data case where the distribution of a binary response is always Bernoulli with mean  $\mu_i$ . The Poisson distribution describes the probability that a specific number of events, say  $y_i$ , occurs,

$$\Pr(Y_i = y_i) = e^{-\mu_i} \mu_i^{y_i} / y_i!; \quad y_i = 0, 1, 2, \dots$$

where  $y! = y \times (y - 1) \times (y - 2) \times \dots \times 2 \times 1$ . The Poisson distribution is completely determined by a single parameter,  $\mu_i = E(Y_i) > 0$ , the mean number of events. A distinctive property of the Poisson distribution is that the mean and variance of  $Y_i$  are equal,

$$E(Y_i) = \mu_i = \text{Var}(Y_i).$$

Note that  $\mu_i$  is defined as the expected count or number of events. The expected rate is given by  $\mu_i/T_i$ , where  $T_i$  is a relevant measure of the “time at risk” (e.g.,  $T_i$  might be an interval of time, the person-years of observation, or the size of a group). In log-linear regression the goal is to describe the effects of a set of covariates,  $X_i$ , on the expected rate. Once again, for ease of exposition, we will first consider the simple case where there is only a single covariate,  $X_i$ . Generalizations to more than one covariate will be considered later.

Because a rate of occurrence of some event cannot be negative, a standard linear regression model relating  $\mu_i/T_i$  directly to  $X_i$  is somewhat unappealing. That is, for sufficiently large or small values of  $X_i$ , a standard linear regression model could yield predicted rates that are negative. Instead, we can relate a transformation of the rate directly to  $X_i$ . When a logarithmic transformation is adopted, the resulting model

$$\log(\mu_i/T_i) = \beta_1 + \beta_2 X_i$$

is known as the *log-linear regression* model. Recall that  $T_i$  is known and fully observed. As a result the log-linear regression model can also be expressed as

$$\log(\mu_i) = \log(T_i) + \beta_1 + \beta_2 X_i,$$

since  $\log(\mu_i/T_i) = \log(\mu_i) - \log(T_i)$ . Note that although  $\log(T_i)$  appears on the right-hand side of the regression equation, it does not have a regression coefficient attached to it. That is, the regression parameter for  $\log(T_i)$  is known to be equal to 1 and does not require estimation. We refer to  $\log(T_i)$  as an *offset*. Thus the log-linear regression model assumes a linear relationship between the log rate of occurrence of some event and  $X_i$ .

When viewed as a generalized linear model, log-linear regression is simply the special case where the distribution of  $Y_i$  is assumed to be Poisson (a member of the exponential family) and a log link

function, the canonical link function for the Poisson distribution, has been adopted. Because the Poisson distribution is a one-parameter exponential family distribution, the variance of  $Y_i$  can be expressed explicitly in terms of the mean, via the following variance function:

$$\text{Var}(Y_i) = v(\mu_i) = \mu_i.$$

For the Poisson distribution, the dispersion parameter is a fixed constant ( $\phi = 1$ ). However, in many biomedical applications, count data have variability that far exceeds that predicted by the Poisson distribution. Overdispersion or extra-Poisson variation is a common indicator of failure of the Poisson assumption when dealing with count data. The implications of overdispersion for count data are the same as for grouped binary data. As discussed earlier, neglecting overdispersion results in the standard errors being underestimated and failure to make an adjustment to the nominal standard errors can result in misleading inferences concerning the regression coefficients (e.g., confidence intervals that are too narrow and  $p$ -values that are too small). Also model selection strategies based on likelihood ratio tests or on information criteria (e.g., AIC) will perform poorly. Adjustments to the nominal standard errors to account for overdispersion can be made either by including a scale factor  $\phi$  in the specification of the Poisson variance,

$$\text{Var}(Y_i) = \phi \mu_i,$$

or by basing standard errors on the “sandwich” estimator of  $\text{Cov}(\hat{\beta})$ ; the “sandwich” estimator will be discussed in greater detail in Chapter 13. A third method of accounting for overdispersion is to incorporate an extra source of variability in the model for the counts; this approach to handling overdispersion is discussed in detail in Section 11.5.

Next consider the interpretation of the log-linear regression coefficients,  $\beta_1$  and  $\beta_2$ . For the special case where the predictor variable  $X_i$  is dichotomous, taking values of 0 and 1, the log-linear regression slope,  $\beta_2$ , has a simple and very attractive interpretation in terms of the log rate ratio (comparing the log expected rate when  $X_i = 1$  to the log expected rate when  $X_i = 0$ ). That is,

$$\log(\mu_i|X_i = 1) - \log(\mu_i|X_i = 0) = \{\log(T_i) + \beta_1 + \beta_2\} - \{\log(T_i) + \beta_1\} = \beta_2.$$

Thus  $\exp(\beta_2)$  has interpretation as the rate ratio

$$\frac{(\mu_i|X_i = 1)}{(\mu_i|X_i = 0)}$$

for the two possible values of the covariate.

For arbitrary  $X_i$  the slope  $\beta_2$  has interpretation as the change in the log expected rate for a single-unit change in  $X_i$ . Equivalently, a unit change in  $X_i$  increases or decreases (depending on the sign of  $\beta_2$ ) the rate of occurrence of the event *multiplicatively* by a factor of  $\exp(\beta_2)$ . Thus, when exponentiated, the regression coefficients can be interpreted in terms of relative rates. This becomes more apparent if we express the log-linear regression model as

$$\mu_i = (\mu_i|X_i) = E(Y_i|X_i) = T_i \times e^{\beta_1} \times e^{\beta_2 X_i}.$$

From this expression it can be seen that a single-unit increase in  $X_i$  increases or decreases  $\mu_i/T_i$  by a factor of  $e^{\beta_2}$ . That is,

$$(\mu_i|X_i + 1) = T_i \times e^{\beta_1} \times e^{\beta_2(X_i + 1)} = T_i \times e^{\beta_1} \times e^{\beta_2 X_i} \times e^{\beta_2} = e^{\beta_2} \times (\mu_i|X_i).$$

On the other hand, the intercept,  $\beta_1$ , has interpretation as the log expected rate when  $X_i = 0$ ; alternatively,  $\exp(\beta_1)$  is the expected rate of occurrence of the event when  $X_i = 0$ .

So far we have only considered the simple case where there is a single predictor  $X_i$ . Next we consider the case where  $X_i$  is a  $p \times 1$  vector of covariates. The log-linear regression model becomes

$$\log(\mu_i) = \log(T_i) + \beta_1 X_{i1} + \beta_2 X_{i2} + \cdots + \beta_p X_{ip},$$

where  $X_{i1} = 1$  for all  $i = 1, \dots, N$ . The log-linear regression coefficients in this model have the following interpretations. Each of the log-linear regression slopes,  $\beta_k$  (for  $k = 2, \dots, p$ ), has interpretation as the change in the log expected rate for a unit change in  $X_{ik}$  given that all of the other covariates remain constant. Thus, holding the remaining covariates at some fixed set of values and not allowing them to vary with any changes in  $X_{ik}$ , we can predict a single-unit increase in  $X_{ik}$  to increase or decrease the log expected rate by an amount  $\beta_k$ . Equivalently, a single-unit increase in  $X_{ik}$

increases or decreases the expected rate *multiplicatively* by a factor of  $\exp(\beta_k)$ . The log-linear regression intercept,  $\beta_1$ , now has interpretation as the log expected rate of occurrence of the event when all covariate values are set to zero. Alternatively,  $\exp(\beta_1)$  is the expected rate when  $X_{i2} = X_{i3} = \dots = X_{ip} = 0$ .

# Illustration

Next we consider an application of log-linear regression to illustrate how the model can be used in practice. The data for this illustration arise from a prospective study of potential risk factors for coronary heart disease (CHD) (Rosenman et al., 1975). The study observed 3154 men aged 40 to 50 for an average of 8 years and recorded the incidence of cases of CHD. The potential risk factors included smoking, blood pressure, and personality/behavior type. The data are summarized in [Table 11.4](#).

**Table 11.4** Data on incidence of CHD and associated risk factors.

Person-Years	Smoking <sup>a</sup>	Blood Pressure <sup>b</sup>	Behavior <sup>c</sup>	CHD
5268.2	0	0	0	20
2542.0	10	0	0	16
1140.7	20	0	0	13
614.6	30	0	0	3
4451.1	0	0	1	41
2243.5	10	0	1	24
1153.6	20	0	1	27
925.0	30	0	1	17
1366.8	0	1	0	8
497.0	10	1	0	9
238.1	20	1	0	3
146.3	30	1	0	7
1251.9	0	1	1	29
640.0	10	1	1	21
374.5	20	1	1	7
338.2	30	1	1	12

*Source:* From *Practical Biostatistical Methods*, 1st edition, by Steve Selvin. © 1995. Reprinted with permission of Brooks/Cole, a division of Thomson Learning: [www.thomsonrights.com](http://www.thomsonrights.com).

<sup>a</sup> 0: Non-smoker, 10: 1–10 cigarettes/day, 20: 11–20 cigarettes/day, 30: 30+ cigarettes/day.

<sup>b</sup> 0: < 140, 1: ≥ 140.

<sup>c</sup> 0: Type B Personality; 1: Type A Personality.

Let  $Y_i$  denote the count of the number of cases of CHD and  $T_i$  denote the person-years of follow-up. Person-years of follow-up is calculated as the total duration of observed follow-up, from entry into the study until either disease detection or end of follow-up, for the individuals in each risk group. To examine whether the rates of CHD are related to the smoking exposure variable we consider the following log-linear regression model:

$$\log(\mu_i/T_i) = \beta_1 + \beta_2 \text{Smoke}_i,$$

where  $\mu_i = E(Y_i)$  and  $\text{Smoke}_i$  is a quantitative measure of smoking exposure (0: Non-smoker, 10: 1–10 cigarettes/day, 20: 11–20 cigarettes/day, 30: 30+ cigarettes/day). To adjust for differences in the total person-years of follow-up for each risk group,  $\log(T_i)$  is included in the model for  $Y_i$  as an offset.

The ML estimate of the slope for smoking exposure,  $\beta_2$ , is 0.0318 (SE = 0.0056) and when compared to its standard error is significantly different from zero at the 0.05 level. This indicates that increases in the smoking exposure increase the log expected rate of CHD. That is, the expected rate of CHD for individuals who smoke one pack of cigarettes (or 20 cigarettes) per day is estimated to be approximately twice (or  $e^{0.0318 \times 20} = 1.88$  times) as high as the rate of CHD for non-smokers.

Because risk factors for CHD are likely to be correlated, we consider the impact of smoking on the rates of CHD after adjusting for the potential confounding effects of blood pressure and personality type. High blood pressure and Type A behavior pattern are known to be associated with high rates of CHD. Specifically, we consider the following log-linear regression model:

$$\log(\mu_i/T_i) = \beta_1 + \beta_2 \text{Smoke}_i + \beta_3 \text{BP}_i + \beta_4 \text{Type}_i,$$

where  $\text{BP}_i = 1$  if blood pressure  $\geq 140$  and 0 otherwise;  $\text{Type}_i = 1$  if Type A personality and  $\text{Type}_i = 0$  if Type B personality.

0 if Type B personality.<sup>1</sup> The estimated log-linear regression parameters (and standard errors) are displayed in [Table 11.5](#).

**Table 11.5** Estimated coefficients and standard errors for log-linear regression of expected rate of CHD on smoking, blood pressure and personality type.

Variable	Estimate	SE	Z
Intercept	-5.4202	0.1308	-41.44
Smoke	0.0273	0.0056	4.88
Personality Type	0.7526	0.1362	5.53
Blood Pressure	0.7534	0.1292	5.83

The estimated coefficient for smoking, 0.027, has now decreased, when blood pressure and personality type have been controlled for in the analysis. Nonetheless, the estimate of  $\beta_2$  remains significantly different from zero at the 0.05 level. The estimated coefficient for smoking has interpretation in terms of the change in the log expected rate of CHD, after adjusting for the effects of blood pressure and personality type. Specifically, the adjusted rate of CHD (controlling for blood pressure and behavior type) for individuals who smoke one pack of cigarettes (or 20 cigarettes) per day is estimated to be 1.7 (or  $e^{0.027 \times 20} = 1.704$ ) times higher than the rate of CHD for non-smokers.

There is also a very strong relationship between Type A behavior pattern and CHD incidence. The adjusted rate ratio (comparing Type A to Type B behavior pattern) is 2.12 (or  $e^{0.7526}$ ), indicating that the rate of CHD among Type A individuals is approximately twice that among Type B individuals. Moreover this adjusted estimate of risk cannot be explained by the association of personality type with smoking and blood pressure since the latter two risk factors have been adjusted for in the analysis.

## 11.4 ORDINAL REGRESSION MODELS

In Section 11.3.1 we discussed regression models for a binary response, a categorical variable with only two levels. In this section we consider regression models for an ordinal response with three or more levels. The reason for devoting a separate section to ordinal regression models is that, when regarded as generalized linear models, they have certain non-standard features. As we will see, ordinal regression models can be regarded as *multivariate*, rather than *univariate*, generalized linear models. Although the main focus of this section is on methods for the analysis of an ordinal response, we conclude this section with a brief discussion of regression models for a nominal or unordered response having more than two levels.

Let us begin by defining what is meant by ordinal data. An ordinal response is a categorical variable whose categories can be naturally ordered, although the precise quantitative distance or spacing between categories is unknown. For example, “cancer stage” is a four-level ordinal variable that describes how far cancer has spread anatomically and can be used to categorize patients with similar prognosis. These four stages are denoted by the roman numerals I through IV, where stage I represents small localized cancers that are usually curable while stage IV represents inoperable or metastatic cancer. A second example of an ordinal variable is socioeconomic status (SES). Socioeconomic status is typically broken into three ordinal categories: high SES, middle SES, and low SES. When defined in this way, a family categorized as “high SES” is higher in the SES hierarchy than a family categorized as “middle SES” (or “low SES”); however, we cannot say “how much higher.” Finally, ordinal variables are frequently used for subjective assessments of quality, importance, or relevance. For example, subjective scales often have response categories of “strongly agree,” “agree,” “disagree,” and “strongly disagree.”

## 11.4.1 Notation

Throughout the remainder of this section we assume that we have  $N$  independent observations of an ordinal response. We let  $Y_i$  ( $i = 1, \dots, N$ ) denote an ordinal response with  $K$  ordinal categories ( $1, \dots, K$ ) for the  $i^{th}$  subject. Note that the actual integer values,  $1, \dots, K$ , are not particularly relevant except that larger values are assumed to correspond to “higher” outcomes and smaller values to “lower” outcomes. The distribution of  $Y_i$  is multinomial (a generalization of the binomial distribution), with  $K$  multinomial probabilities,  $\Pr(Y_i = k)$  for  $k = 1, \dots, K$ , for the distinct ordinal categories; note, there are only  $K - 1$  non-redundant multinomial probabilities because the  $K$  probabilities are constrained to sum to 1. Associated with each response,  $Y_i$ , is a  $p \times 1$  vector of covariates,  $(X_{i1}, \dots, X_{ip})'$ . Next we consider a regression model for the ordinal response that is a direct extension of the familiar logistic regression model for a binary response. However, instead of directly applying the logit transformation to the mean of the ordinal response, the transformation is applied to the *cumulative* response probabilities. This leads to a regression model known as the proportional odds model. The proportional odds model is probably the most widely used model for the analysis of ordinal responses.

## 11.4.2 Proportional Odds Model

To develop the regression model, suppose that we dichotomize the ordinal outcome at 1 versus greater than 1, creating the binary response

$$U_{i1} = \begin{cases} 1 & \text{if } Y_i = 1, \\ 0 & \text{if } Y_i > 1. \end{cases}$$

Letting  $F_{i1} = \Pr(U_{i1} = 1)$ , it is natural to formulate a logistic regression model for the binary response  $U_{i1}$  by relating the logit transformation of  $F_{i1}$  to the covariates,

$$\text{logit}(F_{i1}) = \log\left(\frac{F_{i1}}{1 - F_{i1}}\right) = \alpha_1 + \beta_1 X_{i1} + \beta_2 X_{i2} + \cdots + \beta_p X_{ip};$$

although, in principle, any suitable link function (e.g., probit) could be used. Next, we can dichotomize the ordinal outcome at less than or equal to 2 versus greater than 2, creating a second binary response

$$U_{i2} = \begin{cases} 1 & \text{if } Y_i \leq 2, \\ 0 & \text{if } Y_i > 2, \end{cases}$$

with  $F_{i2} = \Pr(U_{i2} = 1)$ . Because  $U_{i2}$  is also binary, we can formulate a logistic regression model relating the logit of  $F_{i2}$  to the covariates,

$$\text{logit}(F_{i2}) = \log\left(\frac{F_{i2}}{1 - F_{i2}}\right) = \alpha_2 + \beta_1 X_{i1} + \beta_2 X_{i2} + \cdots + \beta_p X_{ip}.$$

Note that we have allowed the intercepts for  $\text{logit}(F_{i1})$  and  $\text{logit}(F_{i2})$  to be different, but have assumed the  $\beta$ 's for the covariates are the same; later we discuss the implications of this assumption. If we continue up the ordinal scale, dichotomizing the ordinal outcome above and below the remaining categories, we can generate a series of additional binary variables

$$U_{ik} = \begin{cases} 1 & \text{if } Y_i \leq k, \\ 0 & \text{if } Y_i > k, \end{cases}$$

and formulate a logistic regression model relating the logit of  $F_{ik}$  to the covariates (for  $k = 1, \dots, K - 1$ ),

$$(11.1) \quad \text{logit}(F_{ik}) = \log\left(\frac{F_{ik}}{1 - F_{ik}}\right) = \alpha_k + \beta_1 X_{i1} + \beta_2 X_{i2} + \cdots + \beta_p X_{ip},$$

where  $F_{ik} = \Pr(U_{ik} = 1 | X_{i1}, \dots, X_{ip}) = \Pr(Y_i \leq k | X_{i1}, \dots, X_{ip})$  is referred to as a “cumulative probability” of response and  $\text{logit}(F_{ik})$  is referred to as a “cumulative log odds” (or “cumulative logit”). The model given by (11.1) is commonly called the *proportional odds model*, and it applies simultaneously to all  $K - 1$  cumulative probabilities (or cumulative logits).

Thus the basic idea underlying the proportional odds model is a cumulative dichotomization of the ordinal variable going up (or down) the ordinal scale. A logistic regression model is assumed to hold simultaneously for each of these  $K - 1$  dichotomous variables, in which the  $K - 1$  intercepts ( $\alpha_k$ 's) are allowed to differ,<sup>2</sup> but the covariate effects ( $\beta$ 's) are assumed to be the same. Therefore the proportional odds model can be thought of as a logistic regression model for the *cumulative probabilities* of response; specifically, it relates the cumulative log odds of response,  $\text{logit}(F_{ik})$ , to the covariates.

Next we consider the interpretation of the regression parameters in the proportional odds model. For ease of exposition, suppose that we have a proportional odds model with two covariates,  $X_{i1}$  and  $X_{i2}$ ,

$$\text{logit}(F_{ik}) = \alpha_k + \beta_1 X_{i1} + \beta_2 X_{i2}.$$

We can interpret  $\beta_1$  in a manner similar to how we interpret regression parameters in standard logistic regression, but recognizing that we are modeling the cumulative log odds (or cumulative logit). Thus  $\beta_1$  has interpretation as the change in the cumulative log odds for each one unit increase in  $X_{i1}$ , while holding  $X_{i2}$  constant. That is, holding  $X_{i2}$  constant,  $\beta_1$  can be thought of as the cumulative log odds ratio,

$$\log\left[\frac{F_{ik}(X_{i1} = c + 1)/\{1 - F_{ik}(X_{i1} = c + 1)\}}{F_{ik}(X_{i1} = c)/\{1 - F_{ik}(X_{i1} = c)\}}\right].$$

The sign of the coefficient for  $\beta_1$  sometimes causes confusion about the direction of the relationship between the ordinal response and the covariate. Recall that (11.1) models the cumulative log odds of being in *lower-numbered* categories. Therefore larger values of  $\beta_1 X_{i1}$  are associated with an *increased* probability of being in the *lower-numbered* categories or, equivalently, a *decreased* probability of being in the *higher-numbered* categories. For example, when  $\beta_1$  is positive this implies an inverse or negative relationship between  $X_{i1}$  and  $Y_i$ , with increases in  $X_{i1}$  associated with lower values of the ordinal scale. We caution the reader that some textbooks, statistical software, and alternative derivations of the model use the following convention for the proportional odds model,

$$\begin{aligned}\text{logit}(F_{ik}) &= \alpha_k - (\beta_1 X_{i1} + \beta_2 X_{i2} + \cdots + \beta_p X_{ip}) \\ &= \alpha_k - \beta_1 X_{i1} - \beta_2 X_{i2} - \cdots - \beta_p X_{ip},\end{aligned}$$

placing a negative sign in front of the  $\beta$ 's so that larger values of  $\beta_1 X_{i1}$  are associated with an *increased* probability of being in the *higher-numbered* categories.

The alert reader would have noticed that our interpretation of  $\beta_1$  is not specific about which of the  $K - 1$  cumulative log odds it refers to. The reason for this is that  $\beta_1$  has the same interpretation for all  $K - 1$  cumulative log odds. In the proportional odds model, the log odds ratio for a one unit increase in a covariate (while holding the other covariates constant) is the same for any of the cumulative probabilities,

$$\begin{aligned}\beta_1 &= \log \left[ \frac{F_{i1}(X_{i1} = c+1)/\{1 - F_{i1}(X_{i1} = c+1)\}}{F_{i1}(X_{i1} = c)/\{1 - F_{i1}(X_{i1} = c)\}} \right] \\ &= \log \left[ \frac{F_{i2}(X_{i1} = c+1)/\{1 - F_{i2}(X_{i1} = c+1)\}}{F_{i2}(X_{i1} = c)/\{1 - F_{i2}(X_{i1} = c)\}} \right] \\ &= \log \left[ \frac{F_{i3}(X_{i1} = c+1)/\{1 - F_{i3}(X_{i1} = c+1)\}}{F_{i3}(X_{i1} = c)/\{1 - F_{i3}(X_{i1} = c)\}} \right].\end{aligned}$$

What this means is that if a unit increase in a covariate triples the odds of being in response level 1 (versus level 2 or higher), it also triples the odds of being in response level 2 or below (versus level 3 or higher), or in level 3 or below (versus level 4 or higher), and so on. This is the “proportionality assumption” that gives the model its name. This property of the proportional odds model also implies that if you were to dichotomize an ordinal response (above and below a given level  $k$ ) and use standard logistic regression as the method of analysis, the resulting odds ratios would be invariant to where you dichotomize the ordinal scale; only the intercept would depend on where you chose to dichotomize the scale.

One appealing property of the proportional odds model is that its regression parameters are invariant to collapsing of adjacent response categories. That is, we would not expect the results of an analysis to change much if we were to combine two adjacent categories. This feature of the model can be helpful when it is of interest to compare regression estimates from studies using different ordinal scales (e.g., one study using a five-level version of a quality-of-life scale, another using a three-level version of the same scale). Also this property of the model can be used to justify combining adjacent categories prior to analysis when data are sparse for certain response categories.

Recall from Section 11.3.1 that the logistic regression model for a binary response can be developed from the notion of an underlying latent variable distribution. The proportional odds model can also be developed from a linear regression model for a latent continuous variable. Suppose that  $L_i$  is a latent (i.e., unobserved) continuous variable, such that values of the ordinal response are observed only when  $L_i$  falls within one of  $K$  intervals determined by a set of “cut-points,”  $\alpha_k$ . In this latent variable formulation, the observed ordinal response can be thought of as a  $K - 1$  level categorization of the unobserved latent variable, with

$$Y_i = \begin{cases} 1 & \text{if } -\infty < L_i \leq \alpha_1, \\ 2 & \text{if } \alpha_1 < L_i \leq \alpha_2, \\ 3 & \text{if } \alpha_2 < L_i \leq \alpha_3, \\ \vdots & \\ K & \text{if } \alpha_{K-1} < L_i < \infty. \end{cases}$$

Next suppose that the following linear regression model holds for  $L_i$ .

$$L_i = \beta_1 X_{i1} + \beta_2 X_{i2} + \cdots + \beta_p X_{ip} + e_i = X'_i \beta + e_i,$$

where  $e_i$  (or  $L_i - X_i\beta$ ) has a standard logistic distribution with mean zero and variance  $\pi^2/3$ . Here we regard the “cut-points” ( $\alpha_k$ ’s) of  $L_i$  as fixed and the mean of the distribution of  $L_i$  as changing with  $X_i$ . It can be shown that this linear model for the latent variable (with logistic errors) implies a proportional odds model for the observed ordinal response with the same covariate effects ( $\beta$ ). Therefore the latent variable formulation provides at least some motivation for the assumption of common covariate effects across the different cumulative logits in the proportional odds model.

As was mentioned earlier, the proportional odds model makes a strong assumption that the covariate effects ( $\beta$ ’s) are invariant to where you dichotomize the ordinal scale. This makes interpretation of the effects of covariates relatively straightforward when only a single parameter is required for each covariate. It is possible to relax the proportionality assumption and consider a “non-proportional” odds model, in which the covariate effects depend on the response level  $k$ ,

$$(11.2) \text{ logit}(F_{ik}) = \alpha_k + \beta_{k1}X_{i1} + \beta_{k2}X_{i2} + \cdots + \beta_{kp}X_{ip};$$

this model allows separate covariate effects for each cumulative logit. In the model given by (11.2) the log odds ratio now depends on  $k$ ,

$$\beta_{k1} = \log \left[ \frac{F_{ik}(X_{i1} = c+1)/\{1 - F_{ik}(X_{i1} = c+1)\}}{F_{ik}(X_{i1} = c)/\{1 - F_{ik}(X_{i1} = c)\}} \right].$$

Model (11.2) can be used to test the proportionality assumption based on a test of the null hypothesis,

$$H_0 : \beta_{1j} = \beta_{2j} = \cdots = \beta_{K-1,j} = \beta_j \text{ for all } p \text{ covariates } (j = 1, \dots, p).$$

Under the null hypothesis that the proportional odds model holds, there are  $p$  distinct  $\beta$ ’s for the covariate effects. Under the alternative hypothesis there are  $(K - 1) \times p$  distinct  $\beta$ ’s. So the test of the proportionality assumption has  $df = (K - 1) \times p - p = (K - 2) \times p$ . Furthermore the proportionality assumption can be relaxed for only a subset of the covariates; this leads to a *partial* proportional odds model where separate effects for each cumulative logit are fit for some but not all of the covariates. For example, the following *partial* proportional odds model,

$$\text{logit}(F_{ik}) = \alpha_k + \beta_{k1}X_{i1} + \beta_2X_{i2} + \beta_3X_{i3} + \cdots + \beta_pX_{ip},$$

allows for separate (or “non-proportional”) effects of  $X_{i1}$  but makes the proportionality assumption for the remaining covariates,  $X_{i2}$ , ...,  $X_{ip}$ . One word of caution about model (11.2) and *partial* proportional odds models: By relaxing the proportionality assumption, the model no longer constrains the cumulative probabilities. As a result the fitting of model (11.2) can potentially lead to incoherent results where, for example, the estimate of  $F_{i3} = \Pr(Y_i \leq 3 | X_{i1}, \dots, X_{ip})$  is less than the estimate of  $F_{i2} = \Pr(Y_i \leq 2 | X_{i1}, \dots, X_{ip})$  for some values of the covariates. This violates the proper order of cumulative probabilities and implies that  $\Pr(Y_i = 3 | X_{i1}, \dots, X_{ip}) = (F_{i3} - F_{i2})$  must be negative!

Finally, the regression parameters of the proportional odds model can be estimated by maximum likelihood (ML). This requires maximizing the multinomial likelihood for the ordinal response, with the response probabilities viewed as functions of the  $\alpha$ ’s and  $\beta$ ’s. When regarded as a generalized linear model, the proportional odds model has certain non-standard features. In the proportional odds model we do not relate the mean of  $Y_i$  to the covariates via a logit link function (or any other suitable link function). Instead, we *jointly* relate the means of the  $K - 1$  cumulative random variables,  $(U_{i1}, \dots, U_{i,K-1})$ , to the covariates. Put another way, it is the cumulative probabilities, and not the mean of the ordinal response, that are simultaneously related to the covariates (via a logit link function). Fitting the proportional odds model requires computer algorithms that can handle the fact that it is a *multivariate*, rather than a *univariate*, generalized linear model for the  $K - 1$  cumulative dichotomizations of the ordinal response. Procedures for fitting the model have been implemented in most of the commercially available statistical software packages.

### 11.4.3 Some Alternative Models for Ordinal Data

Although the proportional odds model is probably the most widely used model for the analysis of ordinal responses, several alternative regression models are used. Here we briefly consider two models, both based on a logistic regression model for the ordinal response, known as the *adjacent-category* and *continuation-ratio* models. The basic idea underlying the *adjacent-category* logistic regression model is to compare each category of the response to the next largest level. So, with a  $K$ -level ordinal response we compare level 1 versus level 2, level 2 versus level 3, level 3 versus level 4, and so on. When these comparisons are made on the log odds scale, the following *adjacent-category* model is obtained,

$$\log \left\{ \frac{\Pr(Y_i = k | Y_i = k \text{ or } k+1)}{\Pr(Y_i = k+1 | Y_i = k \text{ or } k+1)} \right\} = \alpha_k + \beta_1 X_{i1} + \beta_2 X_{i2} + \cdots + \beta_p X_{ip},$$

where the left-hand side is the log odds of response at level  $k$ , given that the response is at either level  $k$  or level  $k+1$ . This model assumes the effects of the covariates do not depend on the particular pair of adjacent categories being compared.

Instead of comparing each category of the response to the next largest level, the *continuation-ratio* model compares each category to all higher response levels. So with a  $K$ -level ordinal response we compare level 1 versus levels 2 through  $K$ , level 2 versus levels 3 through  $K$ , level 3 versus levels 4 through  $K$ , and so on. When these comparisons are made on the log odds scale, the following *continuation-ratio* model is obtained,

$$\begin{aligned} \text{logit}\{\Pr(Y_i = k | Y_i \geq k)\} &= \log \left\{ \frac{\Pr(Y_i = k)}{\Pr(Y_i > k)} \right\} \\ &= \alpha_k + \beta_{k1} X_{i1} + \beta_{k2} X_{i2} + \cdots + \beta_{kp} X_{ip}. \end{aligned}$$

This model ordinarily assumes separate effects of the covariates on each of the  $K-1$  logits; when it seems plausible, it is possible to constrain covariate effects to be the same for each of the logits. The continuation-ratio model can be appealing when the categories of the ordinal response represent a natural sequence of stages in some progression (e.g., cancer stage). One unappealing feature of the model is that the results are not invariant to whether the categories have been ordered from low to high or from high to low.

## 11.4.4 Illustration

We illustrate the application of the proportional odds model to data from a clinical trial comparing auranofin therapy (3 mg of oral gold, twice daily) and placebo for the treatment of rheumatoid arthritis (Bombardier et al., 1986). In this six-month, randomized, double-blind trial, 303 patients with classic or definite rheumatoid arthritis were randomized to one of the two treatment groups and followed over time. The outcome variable of interest is a global impression scale (Arthritis Categorical Scale) at month 6. This is a self-assessment of a patient's current arthritis, measured on a five-level ordinal scale: (1) very good, (2) good, (3) fair, (4) poor, and (5) very poor. Data on this outcome variable are available for 293 of the patients who participated in this trial. The goal of the analysis is to determine whether treatment with auranofin therapy increases the odds of a more favorable response, after controlling for the baseline age of the patients. Consider the following proportional odds model:

$$\text{logit}(F_{ik}) = \log \left\{ \frac{\Pr(Y_i \leq k)}{\Pr(Y_i > k)} \right\} = \alpha_k + \beta_1 \text{Age}_i + \beta_2 \text{Trt}_i$$

where the covariates are baseline age in units of 10 years (Age) and treatment group (Trt = 1 if randomized to auranofin, Trt = 0 if randomized to placebo). Maximum likelihood estimates of the model parameters are presented in [Table 11.6](#).

**Table 11.6** ML estimates and standard errors from the proportional odds model for the arthritis clinical trial data.

Variable	Estimate	SE	Z
$\alpha_1$	-1.2118	0.5316	-2.28
$\alpha_2$	0.5251	0.5249	1.00
$\alpha_3$	2.1494	0.5381	3.99
$\alpha_4$	4.1364	0.6029	6.86
Age	-0.2048	0.0983	-2.08
Trt	0.6079	0.2142	2.84

The results in [Table 11.6](#) indicate that there is a significant treatment effect ( $p < 0.005$ ), with patients in the auranofin therapy group having an increased odds of a lower or more favorable response. Specifically, when adjusted for baseline age, patients in the auranofin therapy group have approximately twice (or  $e^{0.6079} = 1.84$ ) the odds of a self-assessment of arthritis at response level  $k$  or lower (corresponding to a more favorable response) relative to patients in the placebo group. The estimated effect of age indicates that older patients in both treatment groups tend to report less favorable response. For example, a 10-year difference in baseline age decreases the odds of a more favorable response by a factor of 0.82 (or  $e^{-0.205}$ ).

Finally, we can assess the assumption of proportionality by considering a more complex model that allows for separate effects of treatment and age on the four dichotomizations of the ordinal response. The resulting test of proportionality yields a chi-square statistic,  $G^2 = 8.55$ , with 6 degrees of freedom ( $p > 0.20$ ). This test has 6 df because it allows for 6 additional regression parameters, 3 additional parameters for the treatment effect, and 3 additional parameters for the age effect. Because the more complex model does not fit significantly better ( $p > 0.20$ ), the proportional odds assumption appears to hold for these data.

## 11.4.5 Regression Models for Nominal Responses

In this section we briefly discuss regression models when the response is nominal or unordered with more than two levels. Following the notation used in previous sections, we let  $Y_i$  denote a nominal response with  $K$  categories ( $1, \dots, K$ ) for the  $i^{\text{th}}$  subject; the  $K$  unordered categories of the nominal variable are arbitrarily labeled with the integers  $1, \dots, K$ . Associated with each response,  $Y_i$ , is a  $p \times 1$  vector of covariates,  $(X_{i1}, \dots, X_{ip})'$ . To develop a regression model relating  $Y_i$  to the covariates, we consider pairing each response category with a baseline or reference category. For convenience, we choose the last category as the reference category. So with  $K$  categories we compare level 1 versus level  $K$ , level 2 versus level  $K$ , level 3 versus level  $K$ , and so on. When these comparisons are made on the log odds scale, the following *baseline-category* logistic model is obtained,

$$\log \left\{ \frac{\Pr(Y_i = k)}{\Pr(Y_i = K)} \right\} = \alpha_k + \beta_{k1}X_{i1} + \beta_{k2}X_{i2} + \dots + \beta_{kp}X_{ip},$$

for  $k = 1, \dots, K - 1$ . The *baseline-category* logistic model is often referred to as the *multinomial* or *polytomous* logistic regression model.

This model is a very direct extension of standard logistic regression in the sense that it can be formulated as a series of  $K - 1$  logistic regressions for dichotomizations comparing the outcome  $Y_i = k$  versus  $Y_i = K$  (for  $k = 1, \dots, K - 1$ ). For example, consider fitting a logistic regression model for the binary outcome  $Y_i = 1$  versus  $Y_i = K$  (ignoring all observations where  $Y_i = 2, \dots, K - 1$ ); then consider fitting a separate logistic regression model for the binary outcome  $Y_i = 2$  versus  $Y_i = K$  (ignoring all observations where  $Y_i = 1, 3, 4, \dots, K - 1$ ), and so on. Rather than fitting  $K - 1$  separate logistic regressions to the data in this manner, the *baseline-category* logistic model jointly fits the  $K - 1$  logistic regression models.

Interestingly the idea of fitting  $K - 1$  separate logistic regressions to the data happens to very closely approximate the baseline-category multinomial logistic model. That is, the estimates obtained from fitting  $K - 1$  separate logistic regressions will be similar to the corresponding ML estimates obtained by maximizing the multinomial likelihood for the nominal responses. In general, the latter method is preferred as it can yield more precise estimates of the regression parameters, although in practice, the differences may not be too discernible.

For the remainder of the book we do not focus on the analysis of nominal outcomes. The reasons for this decision are two-fold. First, in our experience, nominal responses are not commonly encountered in applications. With the exception of continuous responses, binary, ordinal, and count data are by far the most commonly encountered data types in longitudinal studies and are the main focus of subsequent chapters. Second, as noted earlier, the analysis of nominal data can always be considered a series of  $K - 1$  separate, but otherwise standard, logistic regressions. As a result a comprehensive understanding of extensions of logistic regression models to handle longitudinal binary data provides the basis for a broader understanding of longitudinal methods for analyzing nominal data.

# 11.5 OVERDISPERSION

The Poisson distribution is generally considered to be the benchmark for count data. However, as noted in Section 11.3.2, in many applications count data exhibit far greater variability than is predicted by the Poisson distribution. This overdispersion has important implications for inference. In particular, neglecting overdispersion in regression models for count data results in standard errors being underestimated; failure to make an adjustment to the nominal standard errors can result in misleading inferences concerning the regression coefficients. In this section we describe a method of accounting for overdispersion by incorporating an extra source of variability in the model for the counts. Although the focus of this section is on adjustments to regression models for Poisson count data, the same considerations apply to regression models for binomial counts of the number of successes (or to grouped binary data).

Recall that in the standard log-linear regression model for Poisson count data,

$$\log\{E(Y_i|X_i)\} = \log(\mu_i) = \log(T_i) + \beta_1 X_{i1} + \beta_2 X_{i2} + \cdots + \beta_p X_{ip},$$

the variance of  $Y_i$  is expressed explicitly in terms of the mean,

$$\text{Var}(Y_i) = v(\mu_i) = \mu_i,$$

with the dispersion parameter held fixed at  $\phi = 1$ . Because overdispersion is more the rule than the exception with count data, it is advisable to include a scale factor  $\phi$  in the specification of the variance,

$$\text{Var}(Y_i) = \phi \mu_i.$$

The scale factor  $\phi$  can be estimated by the standard Pearson chi-squared statistic divided by its residual degrees of freedom. The estimated scale parameter affects the standard errors only; specifically, they are multiplied by a factor of  $\sqrt{\phi}$ . This method of adjustment for overdispersion is very simple and somewhat ad hoc; however, in practice, it tends to work well. By including a scale factor  $\phi$ , the variance is assumed to increase *linearly* with increases in the mean.

An alternative way to handle overdispersion is to consider it as an additional source of random variability. That is, overdispersion can be thought to arise due to unmeasured factors that vary among individuals. This suggests extending the log-linear model to incorporate an extra source of variability in the counts,

$$\log\{E(Y_i|X_i; e_i)\} = \log(T_i) + \beta_1 X_{i1} + \beta_2 X_{i2} + \cdots + \beta_p X_{ip} + e_i,$$

where the  $e_i$  are random errors. Conditional on these random errors (and the covariates), it is assumed that the counts have a Poisson distribution. We consider two choices of distributions for the random errors,  $e_i$ . First, assume the errors have a normal distribution, with  $e_i \sim N(0, \sigma_e^2)$ . Then, it can be shown that the mean of  $Y_i$ , when averaged over the distribution of these errors (but conditional on the covariates), is given by,

$$\log\{E(Y_i|X_i)\} = \log(\mu_i) = \log(T_i) + \sigma_e^2/2 + \beta_1 X_{i1} + \beta_2 X_{i2} + \cdots + \beta_p X_{ip}.$$

That is, the model for the mean is the same as the standard log-linear model except for the addition of the constant term  $\sigma_e^2/2$ . Thus, in a model where  $X_{i1} = 1$  for all subjects, only the intercept changes (becoming  $\beta_1 + \sigma_e^2/2$ ); all other regression parameters are unchanged. However, the inclusion of the normally distributed random errors implies that the variance of the counts (when averaged over the distribution of  $e_i$ ) is

$$\text{Var}(Y_i) = \mu_i + (e^{\sigma_e^2} - 1)\mu_i^2,$$

which is larger than the mean,  $\mu_i$ , when  $\sigma_e^2 > 0$ . Therefore the inclusion of this additional source of variation in the log-linear model allows for overdispersion relative to Poisson variation. Of note, the model with normal errors does not have a closed-form likelihood; as a result ML estimation of the model parameters is not entirely straightforward and requires the use of so-called numerical integration techniques that can be computationally demanding. Numerical integration techniques are discussed in Chapter 14.

A similar model that allows for overdispersion relative to Poisson variation can be obtained if, instead of a normal distribution, a gamma distribution is assumed for the errors. Specifically, if a

gamma distribution is assumed for the exponentiated errors,  $\exp(e_i)$ , with mean of 1 and variance denoted by  $\theta$ , then it can be shown that the mean of  $Y_i$ , when averaged over the distribution of these errors (but conditional on the covariates), is given by

$$\log\{E(Y_i|X_i)\} = \log(\mu_i) = \log(T_i) + \beta_1 X_{i1} + \beta_2 X_{i2} + \cdots + \beta_p X_{ip},$$

and the corresponding variance of the counts is

$$\text{Var}(Y_i) = \mu_i + \theta\mu_i^2.$$

Therefore, assuming a gamma distribution for the (exponentiated) errors produces the same model for the mean as in a standard log-linear model, but allows for overdispersion relative to Poisson variability. The nature of the overdispersion is the same under the assumption that the errors have normal and gamma distributions, because  $\theta$  corresponds to  $e^{\sigma_e^2} - 1$  in the two expressions for  $\text{Var}(Y_i)$  given above. In both cases the variance is assumed to increase as a *quadratic* function of the mean, allowing for overdispersion when either  $\theta > 0$  or  $\sigma_e^2 > 0$ . One appealing feature of assuming a gamma distribution for the (exponentiated) errors is that the model has a closed-form likelihood. Specifically, if a gamma distribution is assumed, then the distribution of the counts (when averaged over the distribution of  $e_i$ ) is negative binomial and ML estimation of the model parameters is straightforward. That is, overdispersion can be accounted for by assuming the counts have a negative binomial rather than a Poisson distribution and fitting a log-linear model for the mean via a log link function rather than the canonical link function for the negative binomial (the canonical link for the negative binomial distribution is complicated and somewhat difficult to interpret).

So far we have described two distinct ways to account for overdispersion. The first approach is to include a scale factor in the specification of the variance (allowing for overdispersion when  $\phi > 1$ ); the second approach is to incorporate an additional source of random variability in the model for the counts (allowing for overdispersion when either  $\theta > 0$  or  $\sigma_e^2 > 0$ ). These two approaches make distinct assumptions about the relationship between the mean and variance. For the former, the variance is assumed to increase *linearly* with the mean, whereas for the latter, the variance is assumed to increase as a *quadratic* function of the mean. This implies that the two approaches weight observations differently when fitting the same regression model to the data at hand. In particular, because observations are weighted inversely proportional to their variance, the smallest and largest counts may be weighted somewhat differently by these two approaches. However, in practice, the regression parameter estimates tend to be relatively insensitive to these differences in weights.

Finally, we note that one potential source of overdispersion with count data is when there is a discernibly large number of zero counts. This excess of zeros will necessarily inflate the variability relative to that predicted by the Poisson model. In general, models that include an additional source of random variation (e.g., negative binomial model) allow for a larger probability of both high and low counts in the data. These models will predict, for example, more zeros in the data than a standard Poisson model. However, in certain settings, the introduction of an additional source of random variation may only partially explain the excess number of zero counts. In such cases the large number of zero counts needs to be accounted for using regression techniques that explicitly model the production of zero counts. For example, so-called “zero-inflated Poisson” (ZIP) models have been developed to account for the extra zero counts. These models assume that there are two latent (unobserved) groups: one group can be considered the “always-zero group,” the other the “sometimes-zero group.” For example, in a study of the number of visits per year to a primary care physician, individuals can be thought of as belonging to one of two groups: those who would never visit a primary care physician (the “always-zero group”) versus those who would visit a primary care physician whenever they are sufficiently ill (the “sometimes-zero group”). Observed counts of zero could arise from the former group and from a proportion of the latter group (with the proportion determined by the Poisson probability of a zero count). For example, zero counts of the number of visits per year to a primary care physician can arise from those who would never visit a primary care physician (regardless of their illness status) and from those who would visit but were not ill during the year. A second example arises in a study of the number of children women give birth to. In

such a study there are two latent groups: one group of women who are fertile, another group who are infertile (or whose partners are infertile). That is, zero counts of the number of children can arise from fertile women who, although at risk of pregnancy and bearing children, do not have any children, and from women who are infertile (or whose partners are infertile). ZIP models combine a model for the probability of belonging to the two latent groups with a separate model for the Poisson counts; the latter model for the counts implicitly includes only zero counts from those belonging to the “sometimes-zero group.” A more detailed discussion of ZIP models is beyond the scope of this chapter.

# Illustration

In this section we consider the potential impact of overdispersion on the results of a log-linear regression analysis. We return to the data from the prospective study of potential risk factors for coronary heart disease (CHD) considered in Section 11.3.2. For illustrative purposes we focus on the simple model that examines the unadjusted effect of the smoking exposure variable on rates of CHD. In the standard Poisson log-linear regression model,

$$\log(\mu_i) = \log(T_i) + \beta_1 + \beta_2 \text{Smoke}_i,$$

it is assumed that  $\text{Var}(Y_i) = \mu_i = E(Y_i)$ . We consider three alternative approaches for allowing for overdispersion. First, we include a scale factor  $\phi$  in the specification of the standard Poisson variance,

$$\text{Var}(Y_i) = \phi \mu_i,$$

to be estimated from the data. Next, we extend the log-linear model to incorporate an extra source of variability in the counts,

$$\log\{E(Y_i|e_i)\} = \log(T_i) + \beta_1 + \beta_2 \text{Smoke}_i + e_i,$$

where  $e_i$  are random errors. Two distributions for  $e_i$  are considered: normal and gamma. The inclusion of normally distributed random errors,  $e_i \sim N(0, \sigma_e^2)$ , implies that the variance of the counts is

$$\text{Var}(Y_i) = \mu_i + (e^{\sigma_e^2} - 1)\mu_i^2,$$

whereas the assumption of a gamma distribution for the exponentiated errors,  $\exp(e_i)$ , with mean of 1 and variance  $\theta$ , implies that

$$\text{Var}(Y_i) = \mu_i + \theta\mu_i^2.$$

For the model with gamma errors, ML estimation of the model parameters is straightforward. The model can be fit directly to the counts by assuming they have a negative binomial rather than a Poisson distribution and by assuming a log link function rather than the canonical link function for the negative binomial distribution. In contrast, ML estimation of the model with normal errors requires the use of numerical integration (numerical integration techniques will be discussed in Chapter 14). The results of fitting the four models to the data are summarized in [Table 11.7](#).

**Table 11.7** Estimated coefficients and standard errors for log-linear regression of expected rate of CHD on smoking from (a) standard Poisson model, (b) standard model with overdispersion factor  $\phi$ , (c) Poisson model conditional on normal errors, and (d) Poisson model conditional on gamma errors.

Model	Variable	Estimate	SE	Z
(a) Poisson (fixed $\phi = 1$ )	Intercept	-4.7993	0.0885	-54.22
	Smoke	0.0318	0.0056	5.65
(b) Poisson (unrestricted $\phi$ )	Intercept	-4.7993	0.2415	-19.87
	Smoke	0.0318	0.0153	2.07
$(\hat{\phi} = 7.446)$				
(c) Poisson (normal errors)	Intercept	-4.7069	0.2558	-18.40
	Smoke	0.0282	0.0142	1.99
$(\hat{\sigma}_e^2 = 0.3133)$				
(d) Poisson (gamma errors)	Intercept	-4.5265	0.2536	-17.85
	Smoke	0.0263	0.0141	1.86
$(\hat{\theta} = 0.2942)$				

The ML estimate of the slope for smoking exposure,  $\beta_2$ , from the standard Poisson regression model is 0.0318 (SE = 0.0056) and, when compared to its standard error ( $Z = 5.65, p < 0.0001$ ), is significantly different from zero at the conventional 0.05 level. The results of the analysis that includes an overdispersion factor yield  $\hat{\phi} = 7.45$ . This implies the variability of the counts is approximately  $7\frac{1}{2}$  times larger than that predicted by Poisson variation. The analysis adjusts the

standard errors to account for this degree of overdispersion but does not alter the estimates of the regression coefficients. Specifically, the standard errors are simply multiplied by a factor of 2.7 ( $\sqrt{7.45} = 2.73$ ). When compared to its corrected standard error (SE = 0.0153), the estimate of the slope for smoking exposure remains significantly different from zero at the 0.05 level, although the *p*-value ( $Z = 2.07, p \approx 0.039$ ) is substantially larger than it was before adjustment for overdispersion. The results of the analysis that incorporates an extra source of variability, via the inclusion of normally distributed random errors, yield an estimate of the slope for smoking exposure of 0.0282 (SE = 0.0142). Although the estimated slope for smoking exposure is approximately 10% smaller than that obtained from the standard Poisson regression, the standard errors are  $2\frac{1}{2}$  times larger reflecting the correction made for overdispersion. Finally, the results from the analysis that includes random errors from a gamma distribution (or equivalently, assumes a negative binomial distribution for the counts) yield an estimated slope for smoking exposure of 0.0263 (SE = 0.0141). These results are very similar to those obtained assuming normal errors (albeit the test of slope,  $Z = 1.86, p \approx 0.063$ , is no longer significant at the conventional 0.05 level).

When taken together, the results in [Table 11.7](#) indicate that the expected rate of CHD for individuals who smoke one pack of cigarettes (or 20 cigarettes) per day is estimated to be approximately twice (or  $e^{0.03 \times 20} = 1.82$  times) as high as the rate of CHD for non-smokers. The results also indicate that there is substantial overdispersion in these data. Failure to account for the overdispersion results in the standard errors being underestimated by a factor of approximately 2.5, leading to confidence intervals that are too narrow and *p*-values that are too small.

Finally, we note that the results presented in [Table 11.5](#) for the joint analysis of the effects of smoking exposure, blood pressure, and personality type are far less affected by overdispersion. The degree of overdispersion seen earlier in [Table 11.7](#) for the analysis of smoking exposure is dramatically attenuated when the effects of blood pressure and personality type are included in the regression model. Specifically, the estimated overdispersion factor,  $\hat{\phi}$ , decreases from 7.45 to 1.85. This indicates that these two covariates account for much of the excess variability that was evident in the earlier analysis reported in [Table 11.7](#). That is, much of the overdispersion exhibited in the earlier analysis can be thought of as arising from inherent variation among individuals within each of the smoking exposure groups due to differences in blood pressure and personality type; this between-subject heterogeneity within each of the smoking exposure groups is the major cause of the overdispersion in the counts. When blood pressure and personality type are also included in the analysis,  $\hat{\phi} = 1.85$ . Thus, even when all three factors are controlled, the variability in the counts is almost twice as large as that predicted by Poisson variation. However, when the standard errors in [Table 11.5](#) are corrected by multiplying by a factor of 1.36 (or  $\sqrt{1.85} = 1.36$ ), the overall results are similar. For example, uncorrected for overdispersion, the estimated coefficient for smoking, 0.027 (with 95% confidence interval: 0.016, 0.038), is statistically significant ( $Z = 4.50, p < 0.0001$ ); when corrected for overdispersion, it remains significantly different from zero ( $Z = 3.58, p < 0.0005$ ) although the 95% confidence interval (0.012, 0.042) is slightly wider. In general, we recommend reporting results that are corrected for overdispersion, regardless of its magnitude, because they provide a more realistic estimate of the sampling variability.

# 11.6 COMPUTING: FITTING GENERALIZED LINEAR MODELS USING PROC GENMOD IN SAS

To fit generalized linear models we can use the PROC GENMOD procedure in SAS. The GENMOD procedure fits generalized linear models using maximum likelihood estimation. It includes many of the commonly used exponential family distributions for the response variable and a wide variety of link functions for relating the mean response to the covariates. PROC GENMOD can also be used to fit models to correlated responses using the generalized estimating equations approach. This latter aspect of the procedure will be described in Chapter 13.

For example, to fit a logistic regression model to data from two groups (e.g., treatment or exposure groups), we can use the illustrative SAS commands given in [Table 11.8](#). Similarly, to fit a log-linear regression, with an offset, we can use the illustrative SAS commands given in [Table 11.9](#). To fit a log-linear regression model to overdispersed counts, using a negative binomial rather than a Poisson distribution, we can use the illustrative SAS commands given in [Table 11.10](#). To fit a proportional odds regression model to ordinal data, we can use the illustrative SAS commands given in [Table 11.11](#). Finally, we note that the normal distribution and identity link function are the defaults in PROC GENMOD; this corresponds to the standard linear regression model with normal errors.

**Table 11.8** Illustrative commands for logistic regression using PROC GENMOD in SAS.

---

```
PROC GENMOD DESCENDING;
  CLASS group;
  MODEL y=group / DIST=BINOMIAL LINK=LOGIT;
```

---

**Table 11.9** Illustrative commands for log-linear regression, with an offset, using PROC GENMOD in SAS.

---

```
PROC GENMOD;
  CLASS group;
  MODEL y=group/DIST=POISSON LINK=LOG OFFSET=logtime;
```

---

**Table 11.10** Illustrative commands for log-linear regression assuming negative binomial distribution, using PROC GENMOD in SAS.

---

```
PROC GENMOD;
  CLASS group;
  MODEL y=group / DIST=NEGBIN LINK=LOG OFFSET=logtime;
```

---

**Table 11.11** Illustrative commands for proportional odds regression assuming multinomial distribution and cumulative logit link function, using PROC GENMOD in SAS.

---

```
PROC GENMOD;
  CLASS group;
  MODEL y=group / DIST=MULTINOMIAL LINK=CUMLOGIT;
```

---

Next we present a brief description of the most salient parts of the command syntax used in the four illustrations in [Tables 11.8](#) through [11.11](#).

PROC GENMOD <options>;

This statement calls the procedure GENMOD in SAS. It can include an option for specifying the level of the response variable that is modeled. By default, the lower response level is modeled. For a binary response coded (0,1), it is the probability that  $Y = 0$  that is modeled. Use of the DESCENDING option reverses the default ordering of the response levels, resulting in the highest response level being modeled (i.e., the probability that  $Y = 1$  for binary data that are coded as 0 and 1).

CLASS variables;

The CLASS statement is used to identify all variables that are to be regarded as categorical or factors. By default, this statement will create indicator variables for each factor using a reference group coding, with the last level (where “last” here refers to the level with the largest alphanumeric value) regarded as the reference group. Different sort orders for the CLASS variables can be requested by the ORDER=<option> on the PROC GENMOD statement.

MODEL response = <effects> / <options>;

MODEL events/trials = <effects> / <options>;

The MODEL statement specifies the response variable and the covariate effects. The second form of the MODEL statement, with the events/trials syntax, allows the response to be in the form of a ratio of two variables (e.g., counts of the number of successes and the number of trials) and is used for binomial response data. The covariate effects determine the linear predictor and can include both discrete (defined on the CLASS statement) and quantitative (excluded from the CLASS statement) covariates. By default, PROC GENMOD includes a column of 1's for the intercept in the model.

Two important options need to be included on the MODEL statement. The DIST=*keyword* specifies a built-in response variable distribution, from the exponential family, that is assumed for the model. The LINK=*keyword* specifies the choice of built-in link function relating the mean response to the linear predictor. If the LINK=*keyword* is omitted, the default link function is the canonical link function for the distribution specified on DIST=*keyword*. If both the LINK=<option> and the DIST=<option> are omitted, the default is a normal distribution with an identity link function.

A final option that is often required when modeling count data is an offset. The OFFSET=*variable* specifies a variable to be used as an offset. Note that this variable cannot be a CLASS variable and it should not be included as one of the covariates listed on the MODEL statement.

PROC GENMOD provides many options for handling the dispersion parameter,  $\phi$ , in the exponential family distribution. Recall that for many discrete response distributions (e.g., Bernoulli, binomial, and Poisson), the dispersion parameter is a fixed constant ( $\phi = 1$ ) and not a parameter to be estimated. As discussed earlier, in many applications the data display more variability than is predicted by the variance-mean relationship for the assumed distribution of the response. Neglecting overdispersion (e.g., greater variability than that predicted by the binomial or Poisson distributions) results in standard errors being underestimated. To allow for overdispersion, PROC GENMOD provides options for estimating  $\phi$  and making suitable adjustments to standard errors and test statistics. Strictly speaking, in these cases where  $\phi$  is estimated, rather than assumed to be fixed, we no longer have a legitimate distribution for the response variable and the function that is maximized is referred to as a quasi-likelihood function rather than a likelihood function. Alternatively, for overdispersed counts, a negative binomial rather than a Poisson distribution can be assumed (see [Table 11.10](#)). Finally, an adjustment to the nominal standard errors to account for overdispersion can be made by basing standard errors on the “sandwich” estimator of  $\text{Cov}(\hat{\beta})$ ; the “sandwich” estimator will be described in Chapter 13.

## 11.7 OVERVIEW OF GENERALIZED LINEAR MODELS\*

In this section<sup>†</sup> we present a somewhat more technical and detailed overview of generalized linear models that supplements the material presented in Section 11.2. Generalized linear models are a broad class of regression models suitable for analyzing diverse types of univariate responses (e.g., continuous, binary, counts). As was mentioned in Section 11.2, a generalized linear model for  $Y_i$  has a three-part specification:

1. a distributional assumption,
2. a systematic component, and
3. a link function,

and we consider each of these three components in turn.

# Distributional Assumption

Generalized linear models are an extended family of probability models for a univariate response variable,  $Y_i$ . The family of probability distributions, known as the exponential family, includes the normal distribution for a continuous response, the Bernoulli (or binomial) distribution for a binary response, and the Poisson distribution for counts. The exponential family also includes many other distributions, for example, the gamma, beta, and negative binomial distributions.

Any distribution that belongs to the exponential family can be expressed in the same general form. Before we describe that general form we want to emphasize that our motivation for doing so is three-fold. First, we want to demonstrate that probability distributions for seemingly quite different data types (e.g., continuous, binary, and count data) have much in common as members of the exponential family of distributions. Second, we want to emphasize the importance of the canonical “location” parameter in exponential family distributions; the canonical location parameter is closely related to, but generally not equal to, the mean of the distribution. Third, we want to emphasize that the variance of many exponential family distributions depends on the mean, via a “variance function.” We caution the reader that the material in the remainder of this section is somewhat technical in nature, but we strongly encourage the reader to stay the course.

All distributions that belong to the exponential family can be expressed as follows:

$$(11.3) \quad f(y_i; \theta_i, \phi) = \exp \left[ \{y_i \theta_i - a(\theta_i)\} / \phi + b(y_i, \phi) \right],$$

for some specific functions  $a(\cdot)$  and  $b(\cdot)$ . The specific functions  $a(\cdot)$  and  $b(\cdot)$  associated with an exponential family distribution distinguish one member of the family from another. For example, the normal, Bernoulli, and Poisson distributions can all be expressed in the same form, albeit with different functions  $a(\cdot)$  and  $b(\cdot)$ . This expression for the exponential family has two parameters,  $\theta_i$  and  $\phi$ . The first parameter,  $\theta_i$ , is a location parameter (and is sometimes referred to as the “canonical” location parameter); the second parameter,  $\phi$ , is a scale or dispersion parameter. As these terms imply,  $\theta_i$  is related to the mean of the distribution (but  $\theta_i$  is not necessarily the mean), while  $\phi$  is related to the variance. For many distributions for discrete data,  $\phi$  is not a parameter that requires estimation but is a known constant; for other distributions,  $\phi$  is an unknown parameter. When  $\phi$  is known,  $Y_i$  is said to have a one-parameter exponential family distribution, while when  $\phi$  is unknown, it has a two-parameter exponential family distribution.

While many elegant statistical properties can be derived for distributions that belong to the exponential family, the main concept we want to emphasize in this section is that the exponential family provides some unification for distributions that are commonly assumed for seemingly diverse types of responses variables (e.g., probability distributions for continuous and binary responses).

To fix ideas, we will demonstrate how three of the most commonly encountered distributions in biomedical applications, the normal, Bernoulli, and Poisson distributions, can be expressed in the exponential family form given in (11.3). Recall that the probability density function for the normal distribution (see Section 3.2) is usually written as

$$f(y_i; \mu_i, \sigma^2) = (2\pi\sigma^2)^{-1/2} \exp \left\{ - (y_i - \mu_i)^2 / 2\sigma^2 \right\}.$$

However, it is possible to re-arrange the terms in this expression for the normal density to obtain

$$\begin{aligned} f(y_i; \mu_i, \sigma^2) &= \exp \left\{ -1/2 \log (2\pi\sigma^2) \right\} \exp \left\{ - (y_i - \mu_i)^2 / 2\sigma^2 \right\} \\ &= \exp \left\{ - (y_i^2 - 2y_i\mu_i + \mu_i^2) / 2\sigma^2 - 1/2 \log (2\pi\sigma^2) \right\} \\ &= \exp \left[ \{y_i\mu_i - \mu_i^2/2\} / \sigma^2 - 1/2 \{y_i^2/\sigma^2 + \log (2\pi\sigma^2)\} \right]. \end{aligned}$$

When expressed in this form, the normal distribution is seen to be an exponential family distribution with canonical location parameter,  $\theta_i = \mu_i$ , and scale parameter,  $\phi = \sigma^2$  (with  $v(\mu_i) = 1$ ). Also

$$a(\theta_i) = \mu_i^2/2 = \theta_i^2/2,$$

and

$$\begin{aligned} b(y_i, \phi) &= -1/2 \{y_i^2/\sigma^2 + \log (2\pi\sigma^2)\} \\ &= -1/2 \{y_i^2/\phi + \log (2\pi\phi)\}. \end{aligned}$$

Thus, for the normal distribution the location parameter,  $\theta_i$ , happens to be the mean of the response and the scale parameter happens to be the variance.

Two important exponential family distributions for discrete response data are the Bernoulli and the Poisson distributions. The Bernoulli distribution is ordinarily expressed as

$$f(y_i; \mu_i) = \mu_i^{y_i} (1 - \mu_i)^{1-y_i},$$

where  $\mu_i = E(Y_i) = \Pr(Y_i = 1)$ . At first glance it is not obvious that the Bernoulli distribution also belongs to the exponential family. However, the Bernoulli distribution can also be re-expressed as

$$\begin{aligned} f(y_i; \mu_i) &= \mu_i^{y_i} (1 - \mu_i)^{1-y_i} \\ &= \exp \{y_i \log(\mu_i) + (1 - y_i) \log(1 - \mu_i)\} \\ &= \exp [y_i \log \{\mu_i / (1 - \mu_i)\} + \log(1 - \mu_i)]. \end{aligned}$$

When expressed in this form, the Bernoulli distribution is seen to be a one-parameter exponential family distribution with location parameter,

$$\theta_i = \log \{\mu_i / (1 - \mu_i)\} = \text{logit}(\mu_i),$$

and  $\phi = 1$  is simply a fixed and known constant. Finally, the Poisson distribution is ordinarily expressed as

$$f(y_i; \mu_i) = e^{-\mu_i} \mu_i^{y_i} / y_i!$$

but it too can be re-expressed as

$$f(y_i; \mu_i) = e^{-\mu_i} \mu_i^{y_i} / y_i! = \exp \{y_i \log \mu_i - \mu_i - \log(y_i!)\}.$$

When written in this form, the Poisson distribution is also a one-parameter exponential family distribution with location parameter,

$$\theta_i = \log(\mu_i),$$

and  $\phi = 1$ , a fixed and known constant.

The exponential family unifies many probability distributions for diverse types of response variables. Moreover it is possible to derive some elegant statistical properties for distributions belonging to this family. The two properties that we focus on here are the mean and variance of exponential family distributions. It can be shown (although it requires the use of calculus) that the mean of  $Y_i$  can be expressed as

$$E(Y_i) = \mu_i = \frac{\partial a(\theta_i)}{\partial \theta},$$

where  $\frac{\partial a(\theta_i)}{\partial \theta}$  denotes differentiation of the function  $a(\theta_i)$  with respect to  $\theta$ . For readers unfamiliar with calculus,  $\frac{\partial a(\theta_i)}{\partial \theta}$  can simply be thought of as another known function of  $\theta_i$ . Thus  $\mu_i$ , the mean of  $Y_i$ , is simply a known function of  $\theta_i$ , and vice versa. The second property that we are interested in is the variance of exponential family distributions. The variance of  $Y_i$  can be expressed as

$$\text{Var}(Y_i) = \phi \frac{\partial^2 a(\theta_i)}{\partial \theta^2},$$

where  $\frac{\partial^2 a(\theta_i)}{\partial \theta^2}$  (known in calculus as the second derivative of  $a(\theta_i)$  with respect to  $\theta$ ) is simply another known function of  $\theta_i$ . Thus the variance of  $Y_i$  for distributions belonging to the exponential family can be expressed as the product of  $\phi$ , the dispersion parameter, and some known function of  $\theta_i$ . The latter function is referred to as the “variance function.” However, recall that  $\theta_i$  can be expressed as some known function of the mean,  $\mu_i$  (since earlier we showed that  $\mu_i$  is a known function of  $\theta_i$ ). Because  $\theta_i$  and  $\mu_i$  are functionally related to each other, the variance of  $Y_i$  can be expressed as the product of  $\phi$  and some known function of  $\mu_i$ . When expressed in terms of the mean, the variance function is denoted by  $v(\mu_i)$  and

$$(11.4) \quad \text{Var}(Y_i) = \phi v(\mu_i).$$

Thus, for distributions belonging to the exponential family, the variance of  $Y_i$  can be expressed in terms of a scale or dispersion parameter  $\phi$  and some known function of the mean,  $v(\mu_i)$ . For the normal distribution, the variance of  $Y_i$  is

$$\text{Var}(Y_i) = \sigma^2 = \phi,$$

and  $v(\mu_i) = 1$ . For the Bernoulli distribution,

$$\text{Var}(Y_i) = v(\mu_i) = \mu_i(1 - \mu_i),$$

and  $\phi = 1$ ; while for the Poisson distribution,

$$\text{Var}(Y_i) = v(\mu_i) = \mu_i,$$

and  $\phi = 1$ . For one-parameter exponential family distributions (e.g., Bernoulli and Poisson), the variance of  $Y_i$  is simply a known function of the mean,  $\mu_i$ ; that is, the variance is completely determined by the mean response.

In summary, generalized linear models assume that the response,  $Y_i$ , has a probability distribution that belongs to the “exponential family.” This extended family of distributions includes, among others, the normal, Bernoulli, and Poisson distributions. Some exponential family distributions (e.g., Bernoulli and Poisson) have only a single “location” (or canonical) parameter, and this parameter is related to (but it is not necessarily) the mean of the distribution. For one-parameter exponential family distributions, the variance of  $Y_i$  is a known function of the mean, referred to as the variance function. For two-parameter exponential family distributions (e.g., the normal distribution), there is an additional “scale” parameter, often referred to as a dispersion parameter. In two-parameter exponential family distributions the variance can be expressed as the product of the scale parameter and a variance function, where the latter is a known function of the mean.

# Systematic Component

The systematic component of the generalized linear model specifies that the effects of the covariates,  $X_i$ , on the mean of the distribution of  $Y_i$  can be expressed via the following “linear predictor”:

$$\eta_i = \beta_1 X_{i1} + \beta_2 X_{i2} + \cdots + \beta_p X_{ip}.$$

The linear predictor is simply a linear combination of the unknown vector of regression coefficients,  $\beta = (\beta_1, \dots, \beta_p)'$ , and the vector of covariates,  $X_i$ ,

$$(11.5) \quad \eta_i = \sum_{k=1}^p \beta_k X_{ik}.$$

The term “linear,” as used in this context, means that  $\eta_i$  must be linear in the regression parameters.

We remind the reader that the restriction that  $\eta_i$  be linear in the regression parameters does not preclude relationships between the mean response and covariates that are curvilinear or non-linear. This latter type of non-linearity can be accommodated by taking appropriate transformations of the covariates (e.g.,  $\log(X)$ ) and/or by including a polynomial in  $X$ ). The inclusion of transformed covariates does not violate in any way the requirement that  $\eta_i$  be linear in the regression parameters.

# Link Function

Finally, the formulation of a generalized linear model is completed by specifying the connection between the random and systematic components of the model through a “link function.” The link function describes the relation between  $\mu_i$ , the mean of  $Y_i$ , and the linear predictor,  $\eta_i$ , given by (11.5). Specifically, the link function is some known function  $g(\cdot)$  such that

$$(11.6) \quad g(\mu_i) = \eta_i = \beta_1 X_{i1} + \beta_2 X_{i2} + \cdots + \beta_p X_{ip}.$$

In the case of the standard linear regression model, the random and systematic components are directly related, with

$$E(Y_i) = \mu_i = \eta_i = \beta_1 X_{i1} + \beta_2 X_{i2} + \cdots + \beta_p X_{ip}.$$

When viewed as a generalized linear model, the standard linear regression model adopts an identity link function,  $g(\mu_i) = \mu_i$ .

The primary motivation for considering link functions other than the identity is to ensure that the linear predictor produces predictions of the mean response that are within the allowable range. For example, when analyzing a binary response,  $\mu_i$  has interpretation in terms of the probability of “success.” As a result we must have  $0 < \mu_i < 1$  and the identity link is not appealing since, for sufficiently large or small values of the covariates, it can yield predicted probabilities outside of the range from 0 to 1. It is preferable to use a link function that takes a non-linear transformation of  $\mu_i$ , mapping the range of  $\mu_i$  from  $[0,1]$  onto the unrestricted range  $(-\infty, \infty)$ .

In principle, any function  $g(\cdot)$  can be chosen to link the mean of  $Y_i$  to the linear predictor. However, every distribution that belongs to the exponential family has a special link function called the *canonical* link function. The canonical link function is defined as that function  $g(\cdot)$  such that

$$g(\mu_i) = \theta_i,$$

where  $\theta_i$  is the canonical location parameter (recall that  $\mu_i$  is a known function of  $\theta_i$ , and vice versa). Although there is no a priori reason why the covariate effects should necessarily be additive (or linear) on the particular scale defined by the canonical link function, generalized linear models with canonical link functions produce the most widely used regression models in biomedical applications.

For example, the canonical link function for the normal distribution is the identity link function,

$$g(\mu_i) = \mu_i,$$

and this gives the standard linear regression model,

$$\mu_i = \eta_i = \beta_1 X_{i1} + \beta_2 X_{i2} + \cdots + \beta_p X_{ip}.$$

For counts from a Poisson distribution, where we must have  $\mu_i > 0$ , the canonical link function is the log link function,

$$g(\mu_i) = \log(\mu_i),$$

and this gives the log-linear regression model,

$$\log(\mu_i) = \eta_i = \beta_1 X_{i1} + \beta_2 X_{i2} + \cdots + \beta_p X_{ip}.$$

For the Bernoulli distribution, where  $0 < \mu_i < 1$ , the canonical link function is the logistic or logit link function,

$$g(\mu_i) = \log\{\mu_i / (1 - \mu_i)\}$$

and this gives the logistic regression model

$$\log\{\mu_i / (1 - \mu_i)\} = \eta_i = \beta_1 X_{i1} + \beta_2 X_{i2} + \cdots + \beta_p X_{ip}.$$

We can, however, choose other link functions when they seem appropriate to the application at hand. For example, when  $Y_i$  is Bernoulli, we would generally prefer a link function that transforms the interval  $[0, 1]$  on to the entire real line,  $(-\infty, \infty)$ . The complementary log-log link function,

$$g(\mu_i) = \log\{-\log(1 - \mu_i)\},$$

and the probit link function,

$$g(\mu_i) = \Phi^{-1}(\mu_i),$$

where  $\Phi(\cdot)$  is the standard normal cumulative distribution function, both have this property. In some applications it may be of interest to consider a link function that does not transform the range of  $\mu_i$  on to the entire real line  $(-\infty, \infty)$ . For example, in modeling the prevalence of a disease, and the impact

of risk factors, it may be preferable on scientific grounds to use a log link function since the resulting regression coefficients,  $\beta$ , have interpretation in terms of the log relative risk of disease. The relative risk of disease is simply the ratio of the probability of disease when a risk factor is present to the probability of disease when a risk factor is absent. The relative risk is an index of association that is favored by many empirical researchers. Although, in principle, link functions that do not transform the range of  $\mu_i$  on to  $(-\infty, \infty)$  can be adopted, in practice, this can result in problems with predictions that are out of range (e.g., predicted probabilities less than zero or greater than 1) and problems with convergence of the model-fitting algorithm.

In summary, the link function connects the random and systematic components of the generalized linear model. It relates the mean of  $Y_i$  to the linear predictor and determines the scale on which the additive effects of covariates have an impact on the mean response. Each distribution has a special link function called the canonical link function and adoption of the canonical link function gives rise to many of the widely used regression models in biomedical applications (e.g., linear regression for a normally distributed continuous response, logistic regression for a Bernoulli response, and log-linear regression for Poisson counts). In principle, other link functions can be selected and these link functions bear no relationship to the assumed distribution for the response. Instead, a non-canonical link function may be chosen because additivity of the covariate effects is more appropriate on that scale or because it yields regression coefficients,  $\beta$ , that have somewhat more useful interpretations.

# Estimation

Next we very briefly discuss estimation of the regression coefficients in a generalized linear model. This section is somewhat more technical and can be omitted on a first reading of this chapter. Recall that in the standard linear regression setting we estimate the linear regression coefficients using the method of least squares. The least squares criterion chooses values for the regression coefficients that minimize the sum of squared deviations of the observed  $Y_i$  from their predicted values, denoted by  $\hat{Y}_i = \hat{\mu}_i$ , under the assumed regression model,

$$\mu_i = \eta_i = \sum_{k=1}^p \beta_k X_{ik}.$$

The least squares method yields estimates of the regression coefficients that are also the *maximum likelihood* (ML) estimates when  $Y_i$  is assumed to have a normal distribution with constant variance. We use the more general method of maximum likelihood estimation for estimating the parameters of a generalized linear model.

Recall from Chapter 4 (Section 4.2) that the method of maximum likelihood estimation chooses values of the regression coefficients that are most likely (or most probable) to have generated the observed data. This is achieved by maximizing the *likelihood function* for the data. Construction of the likelihood function requires an assumption about the probability distribution of  $Y_i$ . In generalized linear models the response is assumed to have a distribution belonging to the exponential family of distributions. Assuming independent observations of the response,  $Y_i$ , and the covariates  $X_{i1}, X_{i2}, \dots, X_{ip}$  available on  $N$  individuals, the joint probability of  $(Y_1, \dots, Y_N)$  is the product of the  $N$  probability (density) functions. Thus the likelihood function can be expressed as the product

$$L = \prod_{i=1}^N \exp [\{y_i \theta_i - a(\theta_i)\} / \phi + b(y_i, \phi)].$$

It is this function, or equivalently, the logarithm of this likelihood function, that must be maximized. Note that the likelihood is a function of the unknown regression coefficients,  $\beta$ , since  $\theta_i$  is a known function of the mean,  $\mu_i$ , and

$$\mu_i = g^{-1} \left( \sum_{k=1}^p \beta_k X_{ik} \right),$$

where  $g^{-1}(\cdot)$  denotes the inverse link function; for instance, if  $g(\cdot) = \log(\cdot)$ , then  $g^{-1}(\cdot) = \exp(\cdot)$ . The maximum likelihood estimates of  $\beta$  are obtained by substituting the expression above for  $\mu_i$  into the likelihood function and finding those values of the regression coefficients that produce the largest value for the likelihood function. Ordinarily the likelihood function has only a single maximum.

Instead of maximizing the likelihood, it is usually more convenient to maximize the log-likelihood. We maximize the log-likelihood with respect to  $\beta$  by taking the derivative of the log-likelihood with respect to  $\beta$ , and then finding the values of  $\beta$  that make those derivatives equal to 0. Given

$$l = \log L = \sum_{i=1}^N [\{y_i \theta_i - a(\theta_i)\} / \phi + b(y_i, \phi)],$$

the derivative of the log-likelihood with respect to  $\beta$  can be shown (with the aid of calculus) to be the vector,

$$\partial l / \partial \beta = \sum_{i=1}^N (\partial \theta_i / \partial \beta) (y_i - \mu_i) / \phi.$$

When a canonical link function,  $g(\mu_i) = \theta_i = \eta_i$ , has been assumed,

$$\partial l / \partial \beta = \sum_{i=1}^N X_i (y_i - \mu_i) / \phi.$$

Solving this set of equations,

$$\sum_{i=1}^N X_i (y_i - \mu_i) = 0,$$

yields the ML estimates of  $\beta$ . In general, this requires an iterative procedure that has been implemented in many statistical software packages (e.g., PROC GENMOD in SAS, the `glm` function

in R and S-Plus, and the `g1m` command in Stata). What is quite remarkable about ML estimation for generalized linear models (with canonical link functions) is that it requires the solution to the exact same set of equations, regardless of the type of response variable.

Finally, estimates of the standard errors of the estimated regression coefficients can readily be obtained using the method of maximum likelihood estimation; in addition, likelihood ratio tests can be constructed by comparing nested models. Interestingly the solution to this set of equations,  $\hat{\beta}$ , is consistent for  $\beta$  (i.e., with very high probability,  $\hat{\beta}$  is close to the population regression parameters  $\beta$  for sufficiently large  $N$ ) even if the variance of  $Y_i$  is misspecified; the only requirement is that the model for the mean response (the link function and linear predictor) has been correctly specified. However, when the variance of  $Y_i$  is misspecified, standard errors for components of  $\hat{\beta}$  should be based on the “sandwich” estimator of  $\text{Cov}(\hat{\beta})$ ; the “sandwich” estimator is discussed in Chapter 13.

## 11.8 FURTHER READING

A general overview of logistic regression, Poisson regression, and generalized linear models can be found in Chapter 14 of Neter et al. (1996). The textbooks by Dobson (1990) and Gill (2000) provide excellent introductions to generalized linear models. Hosmer and Lemeshow (2000) provide an accessible and comprehensive description of logistic regression models for binary data. Agresti (2010) provides a comprehensive description of regression models for ordinal data.

# Bibliographic Notes

Generalized linear models were introduced in a seminal paper by Nelder and Wedderburn (1972). McCullagh and Nelder (1989) is the definitive textbook on this topic, providing a comprehensive description of the theory and application of generalized linear models. Firth (1991) presents a concise but remarkably lucid review of generalized linear models; also see Chapter 2 of Fahrmeir and Tutz (2001) and Chapter 5 of McCulloch and Searle (2001). For a comprehensive survey on recent developments in statistical methods for the analysis of ordinal data, see the review article by Liu and Agresti (2005). Cameron and Trivedi (1998) present an overview of regression models for count data, including negative binomial regression models for overdispersed count data; also see Hilbe (2007) for a comprehensive discussion of negative binomial regression models.

## Problems

**11.1** In an experimental study of patients with bladder cancer conducted by the Veterans Administration Cooperative Urological Research Group (Byar and Blackard, 1978; Wei et al., 1989), patients underwent surgery to remove tumors. Following surgery, patients were randomized to either placebo or treatment with thiotepa. Subsequently patients were examined at 18, 24, 30, and 36 months. For this problem set we focus only on the data for month 18. The response variable is binary, indicating whether or not there is a new tumor ( $Y = 1$ , if new tumor;  $Y = 0$ , if no new tumor) at the 18-month visit. The objective of the analysis is to determine the effect of treatment on tumor recurrence by month 18.

The raw data are stored in an external file: `tumor.dat`

Each row of the data set contains the following three variables:

ID Treatment Y

*Note:* The response variable  $Y$  is coded 1 = new tumor, 0 = no new tumor. The categorical variable Treatment is coded 1 = thiotepa, 0 = placebo.

**11.1.1** Assuming a Bernoulli distribution for the recurrence of tumor at month 18, fit the following logistic regression model relating the mean or probability of recurrence ( $\mu_i$ ) to Treatment:

$$\text{logit}(\mu_i) = \beta_1 + \beta_2 \text{Treatment}_i.$$

**11.1.2** What are the interpretations of  $\beta_1$  and  $\beta_2$ ?

**11.1.3** From the results obtained in Problem 11.1.1, what can you conclude about the effect of treatment on tumor recurrence at month 18?

**11.1.4** What is the *estimated* probability of recurrence of a new tumor among those who received placebo?

**11.1.5** What is the *estimated* probability of recurrence of a new tumor among those who received thiotepa?

**11.1.6** Construct a 95% confidence interval for the log odds ratio, comparing thiotepa to placebo.

**11.1.7** Construct a 95% confidence interval for the odds ratio, comparing thiotepa to placebo.

**11.2** In a clinical trial of patients suffering from epileptic seizures (Thall and Vail, 1990), patients were randomized to receive either a placebo or the drug progabide, in addition to standard therapy. A baseline count of the number of epileptic seizures in an 8-week period prior to randomization was obtained. In addition, counts of the number of epileptic seizures in each of four successive 2-week (post-baseline) treatment periods were obtained. For this problem set, we focus only on the data from the last 2-week treatment period. The goal of the analysis is to make a comparison between the two treatment groups in terms of the counts of the number of seizures in the final 2-week period of the study. The question we want to address is whether treatment with progabide is effective in reducing epileptic seizures.

The raw data are stored in an external file: `seizure.dat`

Each row of the data set contains the following four variables:

ID Treatment Age Y

*Note:* The response variable  $Y$  is a count of the number of epileptic seizures in a 2-week interval. The categorical variable Treatment is coded 1 = progabide, 0 = placebo. The variable Age is the age of each patient (in years) at baseline.

**11.2.1** Assuming a Poisson distribution for the counts, fit the following model relating the mean number of seizures ( $\mu_i$ ) to Treatment:

$$\ln(\mu_i) = \beta_1 + \beta_2 \text{Treatment}_i.$$

**11.2.2** What are the interpretations of  $\beta_1$  and  $\beta_2$ ?

**11.2.3** From the results obtained in Problem 11.2.1, what can you conclude about the effect of progabide in reducing the number of epileptic seizures.

**11.2.4** Construct a 95% confidence interval for the rate ratio, comparing progabide to placebo.

**11.2.5** Redo the analysis in Problem 11.2.1, adjusting for the effect of baseline age of the patient:

$$\ln(\mu_i) = \beta_1 + \beta_2 \text{Treatment}_i + \beta_3 \text{Age}_i.$$

**11.2.6** Based on the results of the analysis for Problem 11.2.5, construct a 95% confidence interval for the age-adjusted rate ratio, comparing progabide to placebo.

**11.2.7** Redo the analysis in Problem 11.2.5, allowing for potential overdispersion (i.e., variability greater than that predicted by the Poisson distribution).

**11.2.8** Construct a 95% confidence interval for the age-adjusted rate ratio, comparing progabide to placebo, after taking account of any potential overdispersion.

**11.3** In a study of mental health conducted on a random sample of 40 adult residents of Alachua County, Florida, mental impairment was measured on a four-level ordinal scale with four categories (well, mild symptom formation, moderate symptom formation, impaired); these data are from Chapter 3 (Table 3.3) of Agresti (2010). The goal of the study was to relate mental impairment to several covariates, including an index of life events. The life events (LE) index is a composite measure of the number and severity of important life events that occurred within the past three years (e.g., birth of a child, new job, divorce, death of a family member). The main objective of the analyses is to assess whether the odds of a more favorable mental impairment response is related to the index of life events. Because socioeconomic status (SES) is considered to be a potential confounding variable, it is also of interest to assess the relationship between mental impairment and life events adjusted for SES.

The raw data are stored in an external file: `impairment.dat`

Each row of the data set contains the following four variables:

ID LE SES Y

*Note:* The ordinal response variable  $Y$ , denoting subjects' reported mental impairment, has four categories coded 1=well, 2=mild symptom formation, 3=moderate symptom formation, and 4=impaired. The categorical variable SES is coded 1 = high SES, 0 = low SES. The variable LE is a quantitative measures of the number and severity of important life events.

**11.3.1** Assuming a multinomial distribution for the ordinal response, fit the following proportional odds model relating mental impairment to life events (LE):

$$\log \left\{ \frac{\Pr(Y_i \leq k)}{\Pr(Y_i > k)} \right\} = \alpha_k + \beta_1 \text{LE}_i.$$

**11.3.2** What is the interpretation of the estimate of  $\beta_1$ ?

**11.3.3** Construct a test of the null hypothesis of no effect of life events on the cumulative log odds of response. What conclusions do you draw about the effect of life events on mental impairment?

**11.3.4** Based on the results from Problem 11.3.1, estimate the odds ratio of a more favorable response for subjects with no life events (LE=0) relative to subjects with 6 life events (LE=6).

**11.3.5** The proportional odds model in Problem 11.3.1 makes the assumption of a common effect of life events ( $\beta_1$ ) across the different cumulative logits. Provide a formal or informal assessment of the "proportionality assumption." What do you conclude?

**11.3.6** Redo the analysis in Problem 11.3.1, adjusting for the effect of SES:

$$\log \left\{ \frac{\Pr(Y_i \leq k)}{\Pr(Y_i > k)} \right\} = \alpha_k + \beta_1 LE_i + \beta_2 SES_i.$$

**11.3.7** Based on the results from Problem 11.3.6, what are the interpretations of the estimates of  $\beta_1$  and  $\beta_2$ ?

**11.3.8** Construct a test of the null hypothesis of no effect of life events on the cumulative log odds of response, after adjusting for SES. What conclusions do you draw about the adjusted effect of life events on mental impairment?

**11.3.9** Combine the two adjacent categories, mild and moderate symptom formation, to form a three category ordinal response. With the three category ordinal response, redo the analysis in Problem 11.3.6:

$$\log \left\{ \frac{\Pr(Y_i \leq k)}{\Pr(Y_i > k)} \right\} = \alpha_k + \beta_1 LE_i + \beta_2 SES_i.$$

**11.3.10** Compare and contrast the estimate of  $\beta_1$  obtained from Problem 11.3.9 with the corresponding estimate of  $\beta_1$  obtained from Problem 11.3.6. Does  $\beta_1$  have the same interpretation in the model from Problem 11.3.9 as it does in the model from Problem 11.3.6?

<sup>1</sup> Type A personalities are characterized by impatience, competitiveness, aggressiveness, a sense of time urgency, and tenseness; Type B personalities are the opposite of Type A and exhibit traits such as being easy going, more relaxed about time, not competitive, and not easily angered or agitated.

<sup>2</sup> Note that because this representation of the proportional odds model includes separate parameters ( $\alpha_k$ ) for the  $K - 1$  intercepts, we no longer assume that  $X_{i1} = 1$  for all  $i$ ; this is a slight departure from the notation used in earlier sections and chapters of the book.

<sup>†</sup> This section provides a more technical presentation of generalized linear models and can be omitted without loss of continuity.

# *Chapter 12*

## *Marginal Models: Introduction and Overview*

### **12.1 INTRODUCTION**

In the previous chapter we reviewed generalized linear models for a single response variable. A straightforward application of these models to longitudinal data is not appropriate, owing to the lack of independence among repeated measures obtained on the same individual. There are, however, a number of ways to extend generalized linear models to handle longitudinal data. All of these procedures account for the within-subject correlation among the repeated measures, though they differ in approach. We will see in Chapters 12 through 16 that the method of accounting for the within-subject association has important ramifications for the interpretation of the regression coefficients in models for discrete longitudinal data. For the linear regression models for continuous responses considered in Part II, the interpretation of the regression coefficients is independent of assumptions made about the correlation among the repeated measures. With discrete longitudinal data this is no longer necessarily the case. Instead, different assumptions about the source of the within-subject association can lead to regression coefficients with quite distinct interpretations. The need to distinguish models according to the interpretation of their regression coefficients has led to the use of the terms “marginal models” and “mixed effects models”; the former are often referred to as “population-average models,” the latter as “subject-specific models.” For the former the target of inference is the population, for the latter the target of inference is the individual. In this chapter, we introduce the main features of marginal models for longitudinal data; the meaning of the term “marginal,” as used in this context, will soon be apparent. In Chapter 13, we discuss estimation of marginal models and present three case studies that illustrate the application of marginal models to longitudinal data. Mixed effects models, specifically, generalized linear models with random effects, are the focus of Chapters 14 and 15.

Because the method of accounting for the within-subject association has consequences for the interpretation of the regression model parameters, the choice of method for analyzing discrete longitudinal data cannot be made through any automatic procedure. Rather, the choice must be made on subject-matter grounds. Different models for discrete longitudinal data have somewhat different targets of inference and thereby address subtly different scientific questions. We return to this important issue in Chapter 16.

In this chapter we consider an approach for extending generalized linear models to longitudinal data that leads to a class of regression models that are known as *marginal models*. The term *marginal* in this context indicates that the model for the mean response depends only on the covariates of interest, and not on any random effects or previous responses. That is, the term *marginal* is used to emphasize that the model for the mean response at each occasion does not incorporate dependence on any random effects or previous responses. This is in contrast to *mixed effects models*, where the mean response depends not only on covariates but also on a vector of random effects. Marginal models provide a very natural way of extending generalized linear models to longitudinal data, and they have frequently been applied in the biomedical and health sciences. Marginal models do not require distributional assumptions for the observations, only a regression model for the mean response. That is, marginal models provide a unified method for analyzing diverse types of longitudinal responses, which avoids making assumptions about the distribution of the vector of responses; the method relies solely on assumptions about how the mean response is related to the covariates. The avoidance of distributional assumptions leads to a method of estimation known as *generalized estimating equations* (GEE). The generalized estimating equations

approach is a general method of estimation for marginal models that will be described in detail in Chapter 13.

In our discussion of marginal models in this and the subsequent chapter, the main focus is on discrete response data, for example, binary responses and counts. However, we also point out connections between marginal models for a continuous response and the methods for longitudinal data analysis presented in Part II. In doing so, we can provide some rationale for why the multivariate normal distributional assumption made in Part II often can be relaxed.

## 12.2 MARGINAL MODELS FOR LONGITUDINAL DATA

We begin our discussion of marginal models by introducing some notation similar to that used in Part II. We assume that  $N$  subjects are measured repeatedly over time. We let  $Y_{ij}$  denote the response variable for the  $i^{th}$  subject on the  $j^{th}$  measurement occasion. The response variable can be continuous, binary, ordinal, or a count. The nature of the response variable does have important implications for model specification; however, the notation does not distinguish among the different types of responses.

We do not require that subjects have the same number of repeated measures or that they are measured at a common set of occasions. To accommodate unbalanced data (i.e., repeated measurements that are not obtained at a common set of occasions), we assume that there are  $n_i$  repeated measurements of the response on the  $i^{th}$  subject and that each  $Y_{ij}$  is observed at time  $t_{ij}$ . Both the longitudinal data structure and the notation are the same as that used in Chapter 8; the only difference is that the response variable is no longer assumed to be continuous. The response variables for the  $i^{th}$  subject can be grouped into an  $n_i \times 1$  vector

$$Y_i = \begin{pmatrix} Y_{i1} \\ Y_{i2} \\ \vdots \\ Y_{in_i} \end{pmatrix}, \quad i = 1, \dots, N;$$

where the vectors of responses,  $Y_i$ , are assumed to be independent of one another (but the repeated measures on the same subject are emphatically not assumed to be independent). Associated with each response,  $Y_{ij}$ , there is a  $p \times 1$  vector of covariates

$$X_{ij} = \begin{pmatrix} X_{ij1} \\ X_{ij2} \\ \vdots \\ X_{ijp} \end{pmatrix}, \quad i = 1, \dots, N; \quad j = 1, \dots, n_i.$$

Each individual has a vector of covariates,  $X_{ij}$ , associated with the response at each occasion,  $Y_{ij}$ . Note that  $X_{ij}$  may include covariates whose values do not change throughout the duration of the study and covariates whose values change over time. The former are referred to as time-invariant or between-subject covariates (e.g., gender and fixed experimental treatments), whereas the latter are referred to as time-varying or within-subject covariates (e.g., time since baseline, current smoking status, and environmental exposures). In the former case, the same values of the covariates are replicated in the corresponding rows of  $X_{ij}$ , for  $j = 1, \dots, n_i$ . In the latter case, the values taken by the covariates can vary over time (for at least some individuals) and the values in the corresponding rows of  $X_{ij}$  can be different at each occasion.

We can group the vectors of covariates into an  $n_i \times p$  matrix of covariates

$$X_i = \begin{pmatrix} X'_{i1} \\ X'_{i2} \\ \vdots \\ X'_{in_i} \end{pmatrix} = \begin{pmatrix} X_{i11} & X_{i12} & \cdots & X_{i1p} \\ X_{i21} & X_{i22} & \cdots & X_{i2p} \\ \vdots & \vdots & \ddots & \vdots \\ X_{in_i1} & X_{in_i2} & \cdots & X_{in_ip} \end{pmatrix}, \quad i = 1, \dots, N,$$

where the rows of  $X_i$  correspond to the covariates associated with the responses at the  $n_i$  different measurement occasions, and the columns of  $X_i$  correspond to the  $p$  distinct covariates. So far we have assumed that each subject has a vector of repeated responses, denoted by  $Y_i$ , and associated with each repeated measure there is a vector of  $p$  covariates which can be grouped into a matrix,  $X_i$ .

Marginal models are primarily used to make inferences about population means. As a result marginal models for longitudinal data separately model the mean response and the within-subject association among the repeated responses. In a marginal model the goal is to make inferences about

the former, whereas the latter is regarded as a nuisance characteristic of the data that must be accounted for to make correct inferences about changes in the population mean response.

A marginal model for longitudinal data has the following three-part specification:

1. The conditional expectation or mean of each response,  $E(Y_{ij}|X_{ij}) = \mu_{ij}$ , is assumed to depend on the covariates through a known link function

$$g(\mu_{ij}) = \eta_{ij} = X'_{ij}\beta.$$

2. The conditional variance of each  $Y_{ij}$ , given the covariates, is assumed to depend on the mean according to

$$\text{Var}(Y_{ij}|X_{ij}) = \phi v(\mu_{ij}),$$

where  $v(\mu_{ij})$  is a known “variance function” (i.e., a known function of the mean,  $\mu_{ij}$ ) and  $\phi$  is a scale parameter that may be known or may need to be estimated. For balanced longitudinal designs, a separate scale parameter,  $\phi_j$ , could be estimated at each occasion; alternatively, the scale parameter could depend on the times of measurement, with  $\phi(t_{ij})$  being some parametric function of  $t_{ij}$ .

3. The conditional within-subject association among the vector of repeated responses, given the covariates, is assumed to be a function of an additional set of association parameters,  $\alpha$  (and also depends on the means,  $\mu_{ij}$ ). For example, the components of  $\alpha$  might represent the pairwise correlations or log odds ratios among the repeated responses. The within-subject association among the responses is described in more detail below.

This three-part specification of a marginal model makes the extension of generalized linear models to longitudinal data more transparent. The first two parts of the marginal model correspond to the standard generalized linear model, albeit with no distributional assumptions about the responses (see Section 11.2). It is the third component, the incorporation of the association among the repeated responses from the same individual, that represents the main extension of generalized linear models to longitudinal data. In principle, this three-part specification of a marginal model can be extended by making full distributional assumptions about the vector of responses,  $Y_i$ . However, in Section 12.4 we discuss why assumptions about the joint distribution of  $Y_i$  are not necessary for estimation of the parameters of the marginal model.

As noted above, the first two components of a marginal model specify the mean and variance of  $Y_{ij}$  following the standard generalized linear model formulation described in Chapter 11, the only difference being that we have a common vector-valued link function relating the vector of mean responses to the covariates. The third component recognizes the characteristic lack of independence among longitudinal data by modeling the association among the repeated responses from the same individual. In describing the third component we have been careful to avoid the use of the term *correlation* for two reasons. First, with a continuous response variable, the correlation is a very natural measure of the linear dependence among the repeated responses. Also the correlations are independent of the mean response, in the sense that the correlations are free to vary from  $-1$  to  $1$ , regardless of the values of the vector of mean responses. However, this is not the case with discrete responses. With discrete responses, the correlations are constrained by the mean responses, and vice versa. The most extreme example arises when the response variable is binary. For binary responses the correlations are restricted to ranges that are determined by the means of the responses (or the probabilities of success). For example, in the bivariate case, if  $\mu_1 = E(Y_1) = \Pr(Y_1 = 1) = 0.2$  and  $\mu_2 = E(Y_2) = \Pr(Y_2 = 1) = 0.8$ , then  $\rho_{12} = \text{Corr}(Y_1, Y_2) \leq 0.25$ . That is, the correlation can be no larger than  $0.25$  when the probabilities of success are  $0.2$  and  $0.8$ . As a result, with discrete responses, the correlation is not the most natural measure of within-subject association. Instead, the odds ratio (or the log odds ratio) is a preferable metric for association among pairs of binary responses. Second, for a continuous response that has a multivariate normal distribution, the correlations, along with the variances and the means, completely specify the joint distribution of the vector of longitudinal responses. This is not the case with discrete data. That is, the vector of means and the covariance matrix (the variances and correlations) do not, in general, completely specify the joint distribution of

discrete longitudinal responses. Instead, the joint distribution requires specification of pairwise (e.g., pairwise odds ratios) and higher-order associations among the responses; this feature of discrete data is discussed in greater detail in Section 12.4.

In a certain sense marginal models are a very natural way to extend generalized linear models, developed for the analysis of independent observations, to the setting of correlated longitudinal responses. Marginal models specify a generalized linear model for the longitudinal responses but also include a model for the within-subject association among the responses. A crucial aspect of marginal models is that the mean response and within-subject association are modeled separately. This separation of the modeling of the mean response and the association among responses has important implications for interpretation of the regression parameters in the model for the mean response. In particular, the regression parameters,  $\beta$ , in the marginal model have *population-averaged* interpretations. That is, they describe features of the mean response in the population and how those features relate to covariates. For example, regression parameters in a marginal model might have interpretation in terms of contrasts of the changes in the mean responses in sub-populations (e.g., different treatment or exposure groups or other population strata defined by covariate values). The interpretation of  $\beta$  is not altered in any way by the assumptions made about the nature or magnitude of the within-subject association. We will return to this point in Chapter 16.

Of note, marginal models do not require distributional assumptions for the observations, only a regression model for the mean response. The avoidance of distributional assumptions can be advantageous, since there is no convenient specification of the joint multivariate distribution of  $Y_i$  for marginal models when the responses are discrete. To avoid distributional assumptions for  $Y_i$  we would apply the method of estimation known as *generalized estimating equations* (GEE). The GEE approach provides a convenient alternative to maximum likelihood estimation; the GEE approach for estimating the parameters of marginal models is described in Chapter 13. In Section 12.4 we present a more detailed discussion of how assumptions about the joint distribution of  $Y_i$  are not required for estimation of the marginal model parameters and why it can be advantageous to avoid making distributional assumptions. The material in Section 12.4 is somewhat technical and can be omitted at first reading without loss of continuity.

Finally, we note that there is an implicit assumption in the first component of a marginal model that is often overlooked. Marginal models assume that the conditional mean of the  $j^{th}$  response, given  $X_{i1}, \dots, X_{in_i}$ , depends only on  $X_{ij}$

$$(12.1) \quad E(Y_{ij}|X_i) = E(Y_{ij}|X_{i1}, \dots, X_{in_i}) = E(Y_{ij}|X_{ij}).$$

This assumption implies that given  $X_{ij}$ , there is no dependence of  $Y_{ij}$  on  $X_{ik}$  for  $k \neq j$ . With time-invariant covariates, this assumption poses no difficulties; it necessarily holds since  $X_{ij} = X_{ik}$  for all occasions  $k \neq j$ . Also with time-varying covariates that are fixed by design of the study (e.g., time since baseline, treatment group indicator in a crossover trial), the assumption also holds since values of the covariates at any occasion are determined a priori by study design and in a manner completely unrelated to the longitudinal response. However, when a time-varying covariate varies randomly over time, the assumption made in (12.1) may not hold. For example, the assumption will be violated when the current value of the response, say  $Y_{ij}$ , given the current covariates  $X_{ij}$ , predicts the subsequent value of  $X_{ij+1}$ . This might arise, for example, in a longitudinal observational study designed to assess the effects of physical exercise on reducing blood glucose levels. If study participants with elevated blood glucose levels,  $Y_{ij}$ , at the  $j^{th}$  occasion subsequently increase their amount of physical activity,  $X_{ij+1}$  (while those with normal blood glucose levels continue to maintain their usual level of physical activity), then the assumption made in (12.1) does not hold. As a result somewhat greater care is required when fitting marginal models with time-varying covariates that are not fixed by design of the study. A more detailed discussion of this issue is postponed until Chapter 13 (see Section 13.5).

## **12.3 ILLUSTRATIVE EXAMPLES OF MARGINAL MODELS**

In the previous section we described how marginal models for longitudinal data have a three-part specification in terms of assumptions concerning (1) the mean response at each occasion, (2) the variance of the response at each occasion, and (3) the pairwise within-subject association among the responses. In this section we consider some examples of marginal models using this three-part specification.

# Example 1: Marginal Model for a Continuous Response

The linear regression model for longitudinal data described in Part II is a special case of the marginal model. It is useful to consider its formulation within the framework and terminology of marginal models. By doing so, the extensions to other types of response variables will become more apparent.

Suppose that  $Y_{ij}$  is a continuous response and that it is of interest to relate changes in the mean response over time to the covariates. An example of a marginal model for  $Y_{ij}$  is given by the following three-part specification:

1. The mean of  $Y_{ij}$  is related to the covariates by an identity link function,

$$\mu_{ij} = \eta_{ij} = X'_{ij}\beta.$$

2. The variance of each  $Y_{ij}$ , given the effects of the covariates, is  $\phi$  and does not depend on the mean response. That is,

$$\text{Var}(Y_{ij}|X_{ij}) = \phi v(\mu_{ij}) = \phi,$$

where  $v(\mu_{ij}) = 1$  and  $\phi$  is a scale parameter that needs to be estimated. This model makes the strong, and often unrealistic, assumption that the variance is homogeneous over time. Alternatively, a separate scale parameter,  $\phi_j$ , could be estimated at the  $j^{th}$  occasion if the longitudinal design is balanced on time.

3. The within-subject association among the vector of repeated responses is modeled by assuming a first-order autoregressive correlation pattern

$$\text{Corr}(Y_{ij}, Y_{ik}|X_{ij}, X_{ik}) = \alpha^{|k-j|},$$

where  $0 \leq \alpha \leq 1$ . In this example it is assumed that the within-subject associations do not depend on the means but only on a single correlation parameter,  $\alpha$ . That is,  $\alpha$  is used to model the pairwise correlations among the responses (which are assumed to be approximately equally separated in time).

This illustration of a marginal model for a continuous response is a special case of the linear regression models for longitudinal data considered in Part II. However, marginal models provide a much broader class of models for continuous responses. For example, the means can be related to the covariates by a link function other than the identity or the variances can be allowed to depend on some known function of the means. Also in this illustration the correlations among the components of  $Y_i$  have been specified as a function of the parameter  $\alpha$  via a first-order autoregressive correlation pattern. The correlations can take on many alternative structures, and this example is but one possible structure; other models for the correlation (e.g., unstructured or equicorrelated correlation patterns) can be adopted.

## Example 2: Marginal Model for Counts

Next suppose that  $Y_{ij}$  is a count and we wish to relate changes in the expected count (or expected rate) to the covariates. Counts are often modeled as Poisson random variables, using a log link function and a Poisson variance function. This motivates the following illustration of a marginal model for  $Y_{ij}$ :

1. The mean of  $Y_{ij}$  is related to the covariates through a log link function,

$$\log(\mu_{ij}) = \eta_{ij} = X'_{ij}\beta.$$

2. The variance of each  $Y_{ij}$ , given the effects of the covariates, depends on the mean response,

$$\text{Var}(Y_{ij}|X_{ij}) = \phi \mu_{ij},$$

where  $\phi$  is a time-invariant scale parameter that needs to be estimated.

3. The within-subject association among the vector of repeated responses is assumed to have an unstructured pairwise correlation pattern

$$\text{Corr}(Y_{ij}, Y_{ik}|X_{ij}, X_{ik}) = \alpha_{jk}.$$

Here a balanced longitudinal design is assumed and the vector of parameters  $\alpha$  represents the pairwise correlations among the responses.

The marginal model specified above is a log-linear regression model, with an extra-Poisson variance assumption. The within-subject association is specified in terms of an unstructured pairwise correlation pattern. Of course, other choices for the link and variance functions are possible; similarly other models for the correlation (e.g., first-order autoregressive correlation pattern) are also possible. In this example the extra-Poisson variance assumption allows the variance to be inflated by a factor  $\phi$  (when  $\phi > 1$ ). In many applications count data have variability that far exceeds that predicted by the Poisson distribution; this phenomenon is referred to as *overdispersion*. Indeed, many statisticians believe that overdispersion is the rule, not the exception, when dealing with count data. The excess variability can be accounted for by including the scale factor  $\phi$  in the specification of the variance.

## Example 3: Marginal Model for a Binary Response

Suppose that  $Y_{ij}$  is a binary response, taking values of 0 (denoting “failure”) or 1 (denoting “success”), and it is of interest to relate changes in  $E(Y_{ij}|X_{ij}) = \Pr(Y_{ij} = 1|X_{ij})$  to the covariates. With a binary response the distribution of each  $Y_{ij}$  is Bernoulli, and the probability of success is often modeled using a logit or probit link function. Recall that for a Bernoulli random variable, the variance is a known function of the mean. This motivates the following illustration of a marginal model for  $Y_{ij}$ :

1. The mean of  $Y_{ij}$ , or probability of success, is related to the covariates by a logit link function,

$$\log\left(\frac{\mu_{ij}}{1 - \mu_{ij}}\right) = \eta_{ij} = X'_{ij}\beta.$$

2. The variance of each  $Y_{ij}$ , given the effects of the covariates, depends on the mean response,

$$\text{Var}(Y_{ij}|X_{ij}) = \mu_{ij}(1 - \mu_{ij})$$

and  $\phi = 1$ .

3. The within-subject association among the vector of repeated responses is assumed to have an unstructured pairwise log odds ratio pattern,

$$\log \text{OR}(Y_{ij}, Y_{ik}|X_{ij}, X_{ik}) = \alpha_{jk},$$

where

$$\text{OR}(Y_j, Y_k) = \frac{\Pr(Y_j = 1, Y_k = 1)}{\Pr(Y_j = 1, Y_k = 0)} \cdot \frac{\Pr(Y_j = 0, Y_k = 0)}{\Pr(Y_j = 0, Y_k = 1)}.$$

The marginal model specified above is a logistic regression model, with a Bernoulli variance assumption,  $\text{Var}(Y_{ij}|X_{ij}) = \mu_{ij}(1 - \mu_{ij})$ , and an unstructured within-subject association specified in terms of pairwise log odds ratios rather than pairwise correlations (recall the discussion in Section 12.2 on why the odds ratio is a preferable metric for association among pairs of binary responses).

## Example 4: Marginal Model for an Ordinal Response

Finally, suppose that  $Y_{ij}$  is an ordinal response with  $K$  categories ( $1, \dots, K$ ) and it is of interest to relate changes in the ordinal response to the covariates. To do so, we can specify a marginal model for the *cumulative response probabilities*. For example, a natural extension of the proportional odds model to longitudinal data is

$$\log \left\{ \frac{\Pr(Y_{ij} \leq k | X_{ij})}{\Pr(Y_{ij} > k | X_{ij})} \right\} = \alpha_k + X'_{ij}\beta, \quad (k = 1, \dots, K - 1).$$

In this model, changes in the  $K - 1$  cumulative logits over time are related to the covariates. Although the model includes  $K - 1$  intercepts ( $\alpha_k$ 's), it assumes that the effects of covariates are the same across the  $K - 1$  cumulative logits; this is equivalent to assuming the covariate effects on the cumulative odds are proportional.

Recall from Section 11.4 that the construction of a generalized linear model for the cumulative probabilities requires treating the ordinal response as a set of  $K - 1$  binary variables. Therefore, with repeated measures of an ordinal response,  $Y_{ij}$  can be replaced by  $K - 1$  binary responses,

$$U_{ijk} = \begin{cases} 1 & \text{if } Y_{ij} \leq k, \\ 0 & \text{if } Y_{ij} > k, \end{cases}$$

for  $k = 1, \dots, K - 1$ . That is, the ordinal response at each occasion,  $Y_{ij}$ , is replaced by a vector of binary variables,  $(U_{ij1}, \dots, U_{ij(K-1)})'$  for the  $K - 1$  dichotomizations of the ordinal response. As before, we index subjects by  $i$  and occasions by  $j$ ; however, we now require a third index  $k$  to distinguish the  $K - 1$  dichotomizations of the ordinal response.

Before constructing a marginal model for the ordinal response, we note that the components of  $(U_{ij1}, \dots, U_{ij(K-1)})'$  are correlated, in the sense that  $\text{Corr}(U_{ijk}, U_{ijk'}) \neq 0$  for  $k \neq k'$ . For example,

$$\text{Corr}(U_{ij1}, U_{ij2}) = \frac{F_{ij1} - F_{ij1}F_{ij2}}{\sqrt{F_{ij1}F_{ij2}(1 - F_{ij1})(1 - F_{ij2})}},$$

where  $F_{ij1} = \Pr(U_{ij1} = 1) = \Pr(Y_{ij} \leq 1)$  and  $F_{ij2} = \Pr(U_{ij2} = 1) = \Pr(Y_{ij} \leq 2)$ . This correlation is a direct consequence of the fact that the  $K$  multinomial response probabilities must necessarily sum to 1. Although the components of  $(U_{ij1}, \dots, U_{ij(K-1)})'$  are correlated, the correlations can be expressed in terms of known functions of the cumulative probabilities,  $F_{ijk}$ .

An example of a marginal model for the cumulative response probabilities is given by the following three-part specification:

1. The cumulative probabilities are related to the covariates through a logit link function,

$$\text{logit}\{\Pr(Y_{ij} \leq k | X_{ij})\} = \alpha_k + X'_{ij}\beta, \quad (k = 1, \dots, K - 1).$$

- Letting  $F_{ijk} = \Pr(U_{ijk} = 1 | X_{ij}) = \Pr(Y_{ij} \leq k | X_{ij})$ , this model can be expressed by relating the mean of  $U_{ijk}$  to the covariates through a logit link function,

$$\text{logit}(F_{ijk}) = \alpha_k + X'_{ij}\beta, \quad (k = 1, \dots, K - 1).$$

2. The variance of each  $U_{ijk}$ , given the effects of the covariates, depends only on the mean response,

$$\text{Var}(U_{ijk} | X_{ij}) = F_{ijk}(1 - F_{ijk}).$$

3. In specifying the within-subject association, the correlations among the components of  $(U_{ij1}, \dots, U_{ij(K-1)})'$  at the  $j^{\text{th}}$  occasion are known functions of the means at that occasion,  $F_{ijk}$ . The associations between the components at different occasions are assumed to have an unstructured pairwise pattern.

The marginal model specified above is a proportional odds model for the repeated ordinal response, with a multinomial variance (and covariance) assumption, and an unstructured within-subject association between pairs of repeated measurements.

The four examples of marginal models considered so far are purely illustrative. They demonstrate how the specification of the three components of a marginal model might differ according to the type

of response variable. However, these four examples should not be considered prescriptions for constructing marginal models; in principle, any suitable link function can be chosen and other assumptions about the variances and within-subject associations can be made. The choices for the three components of a marginal model should reflect statistical and subject-matter considerations. In Chapter 13 we present three case studies that illustrate the application of marginal models to longitudinal data.

## 12.4 DISTRIBUTIONAL ASSUMPTIONS FOR MARGINAL MODELS\*

In Section 12.2 a marginal model was defined in terms of a three-part formulation. This formulation highlights how generalized linear models have been extended to handle longitudinal data. In this section<sup>†</sup> we consider making additional distributional assumptions about the vector of responses,  $Y_i$ . Previously we mentioned that specification of the mean vector and the covariance (or the variance and pairwise associations) does not, in most cases, determine the joint distribution of discrete longitudinal data. That is, the three-part marginal model specification does not determine the joint distribution of  $Y_i$ . As a result the method of maximum likelihood cannot be used for estimation of the parameters in the marginal model without further distributional assumptions. This presents two alternative ways to proceed.

The first is to attempt to enrich the formulation of the marginal model so that full distributional assumptions about  $Y_i$  have been made. Then the likelihood can be specified and the method of maximum likelihood can be used for estimation and inference. However, this poses a number of difficulties. First, unlike the multivariate normal distribution for a continuous response, the joint distribution of  $Y_i$  is not usually specified by the mean vector and covariance matrix. That is, with discrete longitudinal data there is no simple analogue of the multivariate normal distribution. Instead, the joint distribution of  $Y_i$  requires specification of the mean vector and pairwise (or two-way) associations, as well as the three-, four-, and higher-way associations among the responses. As the number of responses increases, the number of association parameters proliferates rapidly. This is best exemplified in the case where  $Y_i$  is a vector of binary responses. When the number of repeated measures  $n_i = 10$ , the joint distribution of  $Y_i$  has 1013 (or  $2^{10} - 10 - 1$ ) two-way, three-way, four-way, and higher-way association parameters. This excessive number of within-subject association parameters will often far exceed the number of subjects enrolled in a longitudinal study. As a result specification of the joint distribution for discrete longitudinal data is inherently difficult. In addition, even in cases where it might be possible to specify the joint distribution of  $Y_i$ , the likelihood is often intractable and maximum likelihood estimation is computationally infeasible. Furthermore procedures for ML estimation of marginal models are not currently incorporated in commercially available general-purpose statistical software packages.

The second alternative is to avoid distributional assumptions about  $Y_i$  altogether and specify the marginal model solely in terms of assumptions about the mean response, the variances, and the pairwise (or two-way) within-subject association. This corresponds to the three components in the formulation given in Section 12.2. This alternative approach has the following three advantages. First, it leads to a method for estimation and inference that does not require any distributional assumptions on  $Y_i$ . As a result the empirical researcher does not have to be concerned that the distribution of  $Y_i$  closely approximates some multivariate distribution. Put another way, there may be a gain in robustness because distributional assumptions on  $Y_i$  are not required. Second, it circumvents the need to specify models for the three-way, four-way, and higher-way associations among the responses. Modeling three-way, four-way, and higher-way associations among the responses is conceptually very difficult, and ordinarily requires a relatively large sample size. Third, it leads to a method of estimation known as generalized estimating equations (GEE); the GEE approach is described in detail in Chapter 13. The GEE approach has become an extremely popular method for analyzing longitudinal data, and for good reasons too. It provides a flexible approach for modeling the mean and the pairwise within-subject association structure. It can handle inherently unbalanced designs and missing data with ease. Finally, the GEE approach is computationally straightforward and has been implemented in existing, widely available statistical software. The one potential drawback that must be acknowledged is that avoidance of distributional assumptions will usually result in some loss of efficiency for estimation of  $\beta$  relative to the optimal, but intractable,

likelihood-based estimates. In addition there are some implications for the assumptions made about missing responses; the latter issue will be addressed in Chapters 17 and 18. However, given that the distinct advantages of this alternative approach far outweigh its drawbacks, this is the approach that we emphasize in Chapter 13.

## **12.5 FURTHER READING**

A very accessible description of marginal models, and the generalized estimating equations approach, can be found in Chapter 11 of the textbook by Agresti (2002).

# Bibliographic Notes

There is an extensive statistical literature on likelihood-based marginal models for discrete longitudinal data. Bahadur (1961) proposed a model for the vector of repeated responses expressed in terms of pairwise and higher-order correlations among the responses. Because of the restrictions on the correlations, alternative multinomial models for the joint distribution of the vector of discrete responses have been proposed where the within-subject association is parameterized in terms of other metrics of association. For example, Dale (1984), McCullagh and Nelder (1989), Lipsitz, Laird, and Harrington (1990), Liang, Zeger, and Qaqish (1992), Becker and Balagtas (1993), Molenberghs and Lesaffre (1994), Lang and Agresti (1994), Glonek and McCullagh (1995), and others have proposed full likelihood approaches where the higher-order moments are parameterized in terms of marginal odds ratios. In closely related work, Ekholt (1991) parameterized the association directly in terms of the higher-order marginal probabilities (see also Ekholt, Smith, and McDonald, 1995). An alternative approach parameterizes the within-subject association in terms of conditional associations, leading to so-called “mixed parameter” models (Fitzmaurice and Laird, 1993; Glonek, 1996; Molenberghs and Ritter, 1996).

<sup>†</sup> This section provides a rationale for the use of the generalized estimating equations (GEE) approach for marginal models presented in Chapter 13. The content of this section is somewhat technical and can be omitted at first reading without loss of continuity.

# *Chapter 13*

## ***Marginal Models: Generalized Estimating Equations (GEE)***

### **13.1 INTRODUCTION**

In the previous chapter we considered an approach for extending generalized linear models to longitudinal data that leads to a class of regression models known as *marginal models*. Marginal models have a three-part specification in terms of a regression model for the mean response, supplemented by assumptions concerning the variance of the response at each occasion and the pairwise within-subject association among the responses. In principle, this three-part specification can be extended by making full distributional assumptions about the vector of responses. However, as discussed in Section 12.4, assumptions about the joint distribution of the vector of responses are not necessary for estimation of the marginal model parameters. Indeed, the avoidance of distributional assumptions can be advantageous, since there is no convenient specification of the joint multivariate distribution when the responses are discrete. It also leads to a method of estimation for marginal models known as *generalized estimating equations* (GEE). As we will see, the GEE approach provides a convenient alternative to maximum likelihood (ML) estimation.

The GEE approach for estimating the parameters of marginal models is described in detail in Section 13.2. In Section 13.3, we briefly review some useful residual diagnostics for assessing the fit of marginal models. In Section 13.4, we present three case studies that illustrate the application of marginal models to longitudinal data. Finally, in Section 13.5, we consider estimation and aspects of interpretation of time-varying covariates in marginal models. When a covariate is time-varying, and varies *randomly* over time, subtle issues arise concerning the interpretation and estimation of its effect. Throughout, we adopt the same notation as was used in Chapter 12.

## 13.2 ESTIMATION OF MARGINAL MODELS: GENERALIZED ESTIMATING EQUATIONS

Since there is no convenient specification of the joint multivariate distribution of  $Y_i$  for marginal models when the responses are discrete, we require an alternative to maximum likelihood estimation. The generalized estimating equations (GEE) approach provides that alternative. The GEE approach is based on the concept of “estimating equations” and provides a very general and unified approach for analyzing correlated responses that can be discrete or continuous. The essential idea behind the GEE approach is to generalize and extend the usual likelihood equations for a generalized linear model for a univariate response by incorporating the covariance matrix of the vector of responses,  $Y_i$ . For the case of linear models (i.e., marginal models with an identity link function), the generalized least squares (GLS) estimator of  $\beta$  discussed in Chapter 4 can be considered a special case of the GEE approach. For marginal models with non-linear link functions, this approach is known as “generalized estimating equations”.

Suppose, as in Section 12.2, that the following marginal model has been assumed:

1. The marginal expectation of the response,  $E(Y_{ij}|X_{ij}) = \mu_{ij}$ , depends on the covariates,  $X_{ij}$ , through a known link function

$$g(\mu_{ij}) = \eta_{ij} = X'_{ij}\beta.$$

2. The variance of each  $Y_{ij}$ , given the covariates, depends on the mean according to

$$\text{Var}(Y_{ij}|X_{ij}) = \phi v(\mu_{ij}),$$

where  $v(\mu_{ij})$  is a known “variance function” (i.e., a known function of the mean,  $\mu_{ij}$ ) and  $\phi$  is a scale parameter that may be known or may need to be estimated. For example, when the response is continuous,  $\phi$  is a scale parameter that needs to be estimated. In contrast, with a binary response,  $\phi$  is known and fixed at 1. For count data,  $\phi$  is often estimated from the data at hand to allow for overdispersion relative to Poisson variability.

3. The *pairwise* (or two-way) within-subject association among the vector of repeated responses, given the covariates, is assumed to be a function of the means,  $\mu_{ij}$ , and an additional set of within-subject association parameters,  $\alpha$ . For example, when the vector of parameters  $\alpha$  represents the pairwise correlations among the responses, the covariances among the responses depend on  $\mu_{ij}(X'_{ij}\beta)$ ,  $\phi$ , and  $\alpha$ . That is, given a model for the pairwise correlations, the corresponding covariance matrix can be constructed as the product of standard deviations and correlations

$$V_i = A_i^{\frac{1}{2}} \text{Corr}(Y_i) A_i^{\frac{1}{2}},$$

where  $A_i$  is a diagonal matrix with  $\text{Var}(Y_{ij}|X_{ij}) = \phi v(\mu_{ij})$  along the diagonal (and  $A_{2i}^{\frac{1}{2}}$  is a diagonal matrix with the standard deviations,  $\sqrt{\phi v(\mu_{ij})}$ , along the diagonal), and  $\text{Corr}(Y_i)$  is the correlation matrix (here a function of  $\alpha$ ). In the parlance of the GEE approach,  $V_i$  is known as a “working” covariance matrix to distinguish it from the true underlying covariance among the  $Y_i$ . That is, the term “working” acknowledges our uncertainty about the assumed model for the variances and within-subject associations; unless they have been correctly modeled, our model for the covariance matrix may not be correct.

There are two important features of this specification of a marginal model that are often overlooked. First, recall from Section 12.2 that there is an implicit assumption in the marginal model for the mean response. Marginal models assume that the conditional mean of the  $j^{th}$  response, given  $X_{i1}, \dots, X_{in_i}$ , depends only on  $X_{ij}$ ,

$$E(Y_{ij}|X_i) = E(Y_{ij}|X_{i1}, \dots, X_{in_i}) = E(Y_{ij}|X_{ij}).$$

The implications of this assumption for time-varying covariates are discussed in detail in Section

13.5. Second, the variance of  $Y_{ij}$  at each occasion is specified in terms of a variance function,  $v(\mu_{ij})$ , and a *single* scale parameter  $\phi$ . In principle, a separate scale parameter,  $\phi_j$ , could be estimated at each occasion for balanced designs; alternatively, the scale parameter could depend on the times of measurement, with  $\phi(t_{ij})$  being some parametric function of  $t_{ij}$ . In practice, a limitation of many of the implementations of the GEE approach in widely available software is that they assume the scale parameter  $\phi$  is time-invariant. This restriction on the scale parameter makes these implementations of the GEE approach unappealing for analyzing longitudinal data when the response variable is continuous and the variance of the repeated measurements is not constant over the duration of the study.

Next we provide some motivation for the GEE approach. Recall from Chapter 4 that the generalized least squares (GLS) estimator of  $\beta$  for the linear model minimizes the objective function

$$\sum_{i=1}^N (y_i - X_i\beta)' \Sigma_i^{-1} (y_i - X_i\beta).$$

Using calculus, it can be shown that if a minimum of this function exists it must solve the following equations:

$$\sum_{i=1}^N X_i' \Sigma_i^{-1} (y_i - \mu_i) = 0,$$

where  $\mu_i = \mu_i(\beta) = X_i\beta$ . (Here  $\mu_i(\beta)$  simply denotes that the mean vector,  $\mu_i$ , depends on  $\beta$ .) For the linear model these equations have the following closed-form solution:

$$\hat{\beta} = \left\{ \sum_{i=1}^N (X_i' \Sigma_i^{-1} X_i) \right\}^{-1} \sum_{i=1}^N (X_i' \Sigma_i^{-1} y_i),$$

and  $\hat{\beta}$  is known as the GLS estimator of  $\beta$ . The GEE estimator of  $\beta$  for marginal models (or generalized linear models for longitudinal data) can be thought of as arising from minimizing the following objective function:

$$(13.1) \quad \sum_{i=1}^N \{y_i - \mu_i(\beta)\}' V_i^{-1} \{y_i - \mu_i(\beta)\},$$

with respect to  $\beta$ , where  $V_i$  is treated as known (by ignoring its dependence on  $\beta$  through  $\mu_i$ ) and  $\mu_i$  is the vector of mean responses, with elements

$$\mu_{ij} = \mu_{ij}(\beta) = g^{-1}(X_{ij}'\beta).$$

Using calculus, it can be shown that if a minimum of the function given by (13.1) exists, then it must solve the following *generalized estimating equations*:

$$(13.2) \quad \sum_{i=1}^N D_i' V_i^{-1} (y_i - \mu_i) = 0,$$

where  $V_i$  is the so-called “working” covariance matrix and  $D_i = \partial \mu_i / \partial \beta$  is the gradient or “derivative” matrix (i.e., the matrix containing the derivatives of  $\mu_i$ , with respect to the components of  $\mu$ ). By the “working” covariance matrix we mean that  $V_i$  approximates the true underlying covariance matrix for  $Y_i$ ; that is,  $V_i \neq \text{Cov}(Y_i)$ , recognizing that  $V_i \neq \text{Cov}(Y_i)$  unless the models for the variances and the within-subject associations are correct. As before, we let the true covariance matrix for  $Y_i$  be denoted by  $\Sigma_i$ . The  $n_i \times p$  matrix  $D_i$  can easily be derived using calculus and can be thought of as a matrix that transforms from the original units of  $Y_i$  (and  $\mu_i$ ) to the units of  $g(\mu_{ij})$ . Recall that  $g(\mu_{ij})$  is the scale on which  $\beta$  has interpretation (e.g., the log odds scale rather than the probability scale when the  $Y_{ij}$  are binary and a logit link function has been assumed). The matrix  $D_i$  is only a function of  $\beta$  (since the  $\mu_{ij}$  only depend on  $\beta$ ). For example, when a canonical link is used,  $D_i' = X_i' A_i$ , where  $A_i$  is a diagonal matrix with  $\text{Var}(Y_{ij}|X_{ij}) = \phi v(\mu_{ij})$  along the diagonal. On the other hand,  $V_i$  is a function of  $\beta$ ,  $\phi$ , and  $\alpha$ , since the diagonal elements of  $V_i$  are the variances and the off-diagonal terms are the “working” covariances. That is, the variances depend on the means, and hence  $\beta$ , via the variance function,  $v(\mu_{ij})$  (they also depend on  $\phi$ ); the covariances among the components of  $Y_i$  depend on both  $\beta$  and  $\alpha$ . As a result the generalized estimating equations are functions of both  $\beta$

and  $\alpha$ . For generalized linear models with non-identity link functions, the GEE have no closed-form solution; instead, the solution requires an iterative algorithm.

The generalized estimating equations described above extend in a natural way to marginal models for ordinal responses. Recall from Section 12.3 that when the response is ordinal construction of a marginal model for the cumulative probabilities requires treating each ordinal response as a set of  $K - 1$  binary variables. With repeated measures of an ordinal response,  $Y_{ij}$  is replaced by  $n(K - 1) \times 1$  vector of binary variables, say  $(U_{ij1}, \dots, U_{ij, K-1})'$ , for the  $K - 1$  dichotomizations of the ordinal response (where  $U_{ijk} = 1$  if  $Y_{ij} \leq k$  and  $U_{ijk} = 0$  if  $Y_{ij} > k$ ). Therefore the generalized estimating equations for ordinal responses are based on the following  $(K - 1)n_i \times 1$  vector of binary responses:  $Y_i = (Y_{i1}, Y_{i2}, \dots, Y_{in_i})$  where each  $Y_{ij} = (U_{ij,1}, \dots, U_{ij,K-1})'$  is a  $(K - 1) \times 1$  vector of binary variables; the dimensions of  $\mu_i$ ,  $V_i$ , and  $D_i$  in (13.2) are modified accordingly.

Because the GEE depend on both  $\beta$  and  $\alpha$ , the following iterative two-stage estimation procedure is required:

1. Given current estimates of  $\alpha$  and  $\phi$ ,  $V_i$  is estimated and an updated estimate of  $\beta$  is obtained as the solution to the generalized estimating equations given by (13.2).
2. Given the current estimate of  $\beta$ , updated estimates of  $\alpha$  and  $\phi$  are obtained from the standardized residuals

$$e_{ij} = (Y_{ij} - \hat{\mu}_{ij}) / \sqrt{v(\hat{\mu}_{ij})}.$$

For example,  $\phi$  can be estimated by

$$\hat{\phi} = \frac{\sum_{i=1}^N \sum_{j=1}^{n_i} e_{ij}^2}{\sum_{i=1}^N n_i}.$$

The pairwise association parameters,  $\alpha$ , can be estimated in a similar way, depending on the model for the within-subject association in the third component of the marginal model. For example, in a balanced design, when the association is expressed in terms of unstructured correlations,  $\alpha_{jk} = \text{Corr}(Y_{ij}, Y_{ik})$  can be estimated by

$$\hat{\alpha}_{jk} = \left( \frac{1}{\hat{\phi} N} \right) \sum_{i=1}^N e_{ij} e_{ik}.$$

Finally, in this two-stage estimation procedure, we usually iterate between steps 1 and 2 until convergence has been achieved; starting or initial estimates of  $\beta$  are usually obtained from fitting a generalized linear model assuming independent observations. This algorithm is computationally quite simple, and the GEE approach has been implemented in many statistical software packages (e.g., PROC GENMOD in SAS, the `gee` and `geepack` packages in R, and the `xtgee` command in Stata).

At convergence,  $\hat{\beta}$ , the solution to the generalized estimating equations, has the following properties:

1.  $\hat{\beta}$  is a consistent estimator of  $\beta$ . That is, with very high probability,  $\hat{\beta}$  is close to the population regression parameters  $\beta$  in large samples (i.e., for sufficiently large  $N$ ). Of note,  $\hat{\beta}$  is a consistent estimator of  $\beta$  whether the within-subject associations have been correctly modeled. That is, for  $\hat{\beta}$  to provide a valid estimate of  $\beta$  we only require that the model for the mean response has been correctly specified. This is an important robustness property of  $\hat{\beta}$  that makes the GEE approach very appealing in many applications.
2. In large samples the sampling distribution of  $\hat{\beta}$  is multivariate normal with mean  $\beta$  and

$$\text{Cov}(\hat{\beta}) = B^{-1} M B^{-1},$$

where

$$\begin{aligned} B &= \sum_{i=1}^N D_i' V_i^{-1} D_i, \\ M &= \sum_{i=1}^N D_i' V_i^{-1} \text{Cov}(Y_i) V_i^{-1} D_i. \end{aligned}$$

Note that  $B$  and  $M$  can be estimated by replacing  $\alpha$ ,  $\phi$ , and  $\beta$  by their estimates, and by replacing  $\text{Cov}(Y_i) = \sum_i$  in  $M$  by  $(Y_i - \hat{\mu}_i)(Y_i - \hat{\mu}_i)'$ . That is, the expression for  $\widehat{\text{Cov}}(\hat{\beta})$  is given by

$$(13.3) \quad \left( \sum_{i=1}^N \hat{D}_i \hat{V}_i^{-1} \hat{D}_i \right)^{-1} \left\{ \sum_{i=1}^N \hat{D}_i \hat{V}_i^{-1} (Y_i - \hat{\mu}_i)(Y_i - \hat{\mu}_i)' \hat{V}_i^{-1} \hat{D}_i \right\} \left( \sum_{i=1}^N \hat{D}_i \hat{V}_i^{-1} \hat{D}_i \right)^{-1}.$$

This is known as the empirical or “sandwich” estimator; the components  $B$  and  $M$  can be thought of as the “bread” and “meat” of this sandwich estimator of  $\text{Cov}(\hat{\beta})$ . Finally, if we model  $V_i$  correctly,  $V_i = \sum_i$ , and  $\text{Cov}(\hat{\beta}) = B^{-1}$ .

In summary, the GEE approach has a number of appealing properties for estimation of the regression parameters in marginal models. First, in many longitudinal designs the GEE estimator  $\hat{\beta}$  is almost as precise or efficient as the MLE. For example, it can be shown that the GEE has a similar expression to the likelihood equations for  $\beta$  in a linear model for continuous responses that are assumed to have a multivariate normal distribution. That is, the GLS estimator of  $\beta$  can be considered a special case of the GEE approach. The GEE also has an expression similar to the likelihood equations for  $\beta$  in certain models for discrete longitudinal data. As a result, for many longitudinal designs, there is little loss of precision when the GEE approach is adopted as an alternative to maximum likelihood. Second, the GEE estimator  $\hat{\beta}$  is a consistent estimator of  $\beta$  even if the within-subject associations among the repeated measures have been misspecified; this is a very appealing robustness property of the GEE estimator. Although the GEE estimator  $\hat{\beta}$  is a consistent estimator under misspecification of the within-subject associations, the usual standard errors obtained under the misspecified model for the within-subject association are not valid. Fortunately, in many cases valid standard errors for  $\hat{\beta}$  can be obtained using the empirical or “sandwich” estimator of  $\text{Cov}(\hat{\beta})$ . In addition to correcting for misspecification of the within-subject association, the “sandwich” estimator also corrects for potential overdispersion (or, indeed, any misspecification of the variance). Another important feature of the GEE approach is that it can readily handle imbalance due to missing data in the response variables. However, in doing so, it does require a strong, and often unrealistic, assumption that data are missing completely at random (MCAR). The GEE estimators can be adapted to provide a valid analysis when data are missing at random (MAR), but not MCAR, by explicitly modeling the missingness process and weighting the analysis accordingly. This approach to handling missing data is known as the inverse probability weighted (IPW) GEE method and is discussed in detail in Chapters 17 and 18.

Finally, although the main emphasis of this chapter has been on longitudinal analysis of a discrete response, the GEE approach can be applied equally to continuous responses. That is, for the linear regression models described in Part II, the multivariate normal assumption is not crucial. Specifying a linear regression model for the longitudinal responses and a model for the covariance among the responses is sufficient for the purposes of estimating  $\beta$  using the GEE approach. As was mentioned above, the GLS estimator of  $\beta$  can be considered a special case of the GEE approach and the multivariate normal distribution assumption for the responses is not required. The validity of the GLS/GEE estimates of  $\beta$  rests only on having a correct model for the mean response. When the model for the covariance is misspecified, valid standard errors for  $\hat{\beta}$  can be obtained using the “sandwich” estimator of  $\text{Cov}(\hat{\beta})$ . Thus, although the GEE approach and the “sandwich” estimator of  $\text{Cov}(\hat{\beta})$  are more widely used in marginal models for discrete data, they can also be applied in the linear models for continuous data described in Part II. However, because many of the implementations of the GEE approach in widely available software assume that the scale parameter  $\phi$  is time-invariant, we do not recommend their use for analyzing longitudinal data when the response variable is continuous. Instead, the GEE approach can be implemented using existing software for the general linear model that allows a much wider range of covariance pattern models and/or random effects covariance structures, coupled with the option of calculating standard errors for  $\hat{\beta}$  based on the “sandwich” estimator (e.g., using the EMPIRICAL option in PROC MIXED in SAS).

# A Note on the “Sandwich” Estimator of $\text{Cov}(\hat{\beta})$

An appealing property of the GEE estimator  $\hat{\beta}$  is that it is a consistent estimator of  $\beta$  even if the assumed model for the covariances among the repeated measures is not correct. It only requires that the model for the mean response be correct. This robustness property of GEE is important because the usual focus of a longitudinal study is on changes in the mean response. Based either on theoretical grounds (e.g., randomization in an experiment) or subject-matter knowledge of similar data, the data analyst can often specify how changes in the mean response depend on the covariates. On the other hand, much less is usually known about the patterns of two-way and higher-way associations among the responses; moreover these “higher-order moments” are increasingly difficult to estimate from the data.

For inferences about  $\beta$ , valid standard errors can be obtained from the “sandwich” estimator of  $\text{Cov}(\hat{\beta})$  given by (13.3). The remarkable property of the “sandwich” estimator is that it is also robust in the sense that it provides valid standard errors when the assumed model for the covariances among the repeated measures is not correct. That is, with large sample sizes, the “sandwich” estimator of  $\text{Cov}(\hat{\beta})$  yields correct standard errors. Because of this appealing robustness property, the “sandwich” estimator is often referred to as the “robust” variance estimator or the “empirical” variance estimator. Maintaining the culinary theme of this section, it would seem that we can have our cake and eat it: we can obtain a valid estimate of  $\beta$  and its sampling variability, even if we have not modeled the within-subject associations correctly. Indeed, some readers may see it as a delicious irony that we can disregard the model for the covariances among the repeated measures for the purposes of inference about  $\beta$ .

This raises an important issue. Why bother expending effort to model the within-subject association? For example, naively assuming the responses are independent (i.e., specifying the “working” covariance matrix as diagonal) yields valid estimates of  $\beta$ ; valid standard errors can then be obtained using the “sandwich” estimator of  $\text{Cov}(\hat{\beta})$ . There are two main reasons for modeling the covariance. First, in general, the closer the “working” covariance matrix ( $V_i$ ) approximates the true underlying covariance matrix ( $\Sigma_i$ ), the greater the efficiency or precision with which  $\beta$  can be estimated. That is, a “working” covariance matrix that approximates the true underlying covariance matrix makes optimal use of the available data for estimation of  $\beta$ . Second, the robustness property of the “sandwich” estimator of  $\text{Cov}(\hat{\beta})$  is a large sample (or asymptotic) property. In general, use of the “sandwich” estimator of  $\text{Cov}(\hat{\beta})$  is best suited to balanced longitudinal designs where the number of subjects ( $N$ ) is relatively large and the number of repeated measures ( $n$ ) is relatively small. Moreover the “sandwich” estimator is less appealing when the design is severely unbalanced and/or when there are few replications to estimate the true underlying covariance matrix. In applications, use of the “sandwich” estimator implicitly relies on there being many replications of the vector of responses associated with each distinct set of covariate values. For example, in a longitudinal clinical trial with two treatment groups there will be many replications of  $Y_i$  associated with the two distinct sets of covariate values ( $X_i$ ) for the treatment and control groups. In that case, use of the “sandwich” estimator of  $\text{Cov}(\hat{\beta})$  is justified because there is a sufficient number of replications (or number of subjects) to estimate the true underlying covariance matrix within each treatment group. In many observational studies, however, there may be few, if any, replications of  $Y_i$  associated with each distinct set of covariate values, especially when  $X_i$  includes many covariates and/or quantitative covariates. Similarly, if the longitudinal design is severely unbalanced, with each individual having a unique sequence of measurement occasions,  $t_{i1}, \dots, t_{in_i}$ , there are no replications at each of the measurement occasions. In these cases the use of the “sandwich” estimator is less appealing. In particular, when the number of subjects is relatively small the “sandwich” -based standard errors are biased downward, in the sense that the nominal standard errors are too small and underestimate the variance of  $\hat{\beta}$ . Moreover the magnitude of the bias of the “sandwich” variance estimator depends on the covariate design matrix,  $X_i$ . In addition, the sampling variability of the “sandwich” estimator of

$\text{Cov}(\hat{\beta})$  can be very large, resulting in an unstable estimate of  $\text{Cov}(\hat{\beta})$ . This increased variability of the “sandwich” estimator is the price paid for its robustness property. However, the increased variability can lead to problems with the coverage probabilities of confidence intervals constructed from “sandwich” estimates. In particular, use of the “sandwich” estimates can yield confidence intervals with coverage probability well below the desired nominal level.

In summary, the “sandwich” estimator of  $\text{Cov}(\hat{\beta})$  is of most practical use when the sample size is relatively large or when the assumed model for the covariances among the repeated measures (the “working” covariance matrix) is questionable. Reliance on the “sandwich” estimator is less appealing when the number of independent subjects is modest (relative to the number of repeated measures), when the design is inherently unbalanced, or when subjects cannot be grouped on the basis of having identical covariate design matrices. For any of these cases it can be advantageous to model the covariances among the responses and use a “model-based” estimator of  $\text{Cov}(\hat{\beta})$ . The model-based estimator is given by

$$\text{Cov}(\hat{\beta}) = B^{-1},$$

where

$$B = \sum_{i=1}^N D_i' V_i^{-1} D_i,$$

and can be estimated by replacing  $\alpha$ ,  $\phi$ , and  $\beta$  by their estimates. This estimator of  $\text{Cov}(\hat{\beta})$  is called a “model-based” estimator to remind us that it yields valid standard errors provided that the “working” covariance matrix,  $V_i$ , is a close approximation to the true underlying covariance matrix,  $\Sigma_i$ . That is, the “model-based” estimator does require that the model for the covariance, the “working” covariance, be correctly specified.

## 13.3 RESIDUAL ANALYSES AND DIAGNOSTICS

The adequacy of a fitted marginal regression model can be assessed using various residual analysis techniques similar to those described in Chapter 10 for linear models for longitudinal data. Residuals and other standard regression diagnostics (e.g., Cook's distance and leverage) can also be helpful for identifying outliers and influential observations and/or influential individuals. In this section we briefly review some useful diagnostics for assessing the functional form of the marginal model for the mean response and for detecting observations or individuals that may be having an undue influence on the analysis.

For a marginal regression model it is straightforward to calculate residuals based on the difference between the observed and predicted responses at each occasion,

$$r_{ij} = Y_{ij} - \hat{\mu}_{ij} = Y_{ij} - g^{-1}(X'_{ij}\hat{\beta});$$

these can be readily produced by most statistical software packages for fitting marginal models. However, because the variance of the residual,  $r_{ij}$ , is a function of the predicted mean response, it is preferable to conduct all model checking using studentized Pearson residuals,

$$e_{ij} = \frac{Y_{ij} - g^{-1}(X'_{ij}\hat{\beta})}{\sqrt{\phi v(\hat{\mu}_{ij})(1 - h_{ij})}},$$

where  $v(\mu_{ij})$  is the variance function,  $\phi$  is the scale parameter (either fixed or estimated from the data, depending on the type of response and assumptions about overdispersion), and  $h_{ij}$  is the “leverage” of the  $j^{th}$  observation on the  $i^{th}$  individual. In regression, the leverage,  $h_{ij}$ , describes the potential influence each observation has on its own predicted value. In general, observations that are extreme in terms of the  $X_{ij}$ 's have high leverage (albeit in generalized linear models, leverage values also depend on the mean response). These Pearson residuals can be used to check for any systematic departures from the model for the mean response; for example, a scatterplot of  $e_{ij}$  against the predicted mean response,  $\hat{\mu}_{ij}$ , can be examined for the appearance of systematic trend. The fitting of a smooth curve (e.g., a *lowess* curve) to the scatterplot can often help in judging whether curvature is present. Similarly scatterplots of the residuals against selected covariates from the model for the mean response can be examined for any systematic trends. Such a trend may indicate, for example, the omission of a quadratic term or the need for transformation of the covariate. A scatterplot of the residuals versus time (or age) can be particularly useful for assessing the adequacy of the marginal model assumptions about patterns of change in the mean response over time.

An acknowledged difficulty with the interpretation of these conventional residual diagnostic plots is that they are inherently subjective. In Chapter 10 (Section 10.4) we discussed how this problem can be overcome by basing model assessment on “cumulative sums” and “moving sums” of residuals. The exact same approach can be extended to residuals from marginal regression models. That is, we can compare the *observed* sum of the residuals, both graphically and numerically, to a reference distribution under the assumption of a correctly specified marginal model for the mean response. This allows us to determine whether any apparent pattern in the observed sum of the residuals is evidence of a systematic trend or simply due to natural variation. For example, to check the functional form of the  $k^{th}$  covariate, we can define the cumulative sum of the residuals over values of  $X_{ijk}$ ,

$$W_k(x) = \frac{1}{\sqrt{N}} \sum_{i=1}^N \sum_{j=1}^{n_i} I(X_{ijk} \leq x) r_{ij},$$

where  $I(\cdot)$  is the indicator function. In addition we can construct the cumulative sum of residuals over the fitted values, denoted by  $W_f(x)$ ,

$$W_f(x) = \frac{1}{\sqrt{N}} \sum_{i=1}^N \sum_{j=1}^{n_i} I(X'_{ij}\hat{\beta} \leq x) r_{ij}.$$

Recall from Section 10.4 that if the assumed model for the marginal mean has been correctly

specified, the cumulative residual process should be centered at zero and behave like a zero-mean Gaussian (or normal) process. This zero-mean Gaussian process provides a reference for deciding whether any pattern in the observed cumulative residual process is systematic or simply due to the natural variation of the process. An assessment of model adequacy is based on comparing the pattern of the observed cumulative residual process with computer simulated realizations from the zero-mean Gaussian process (the null distribution). In particular, the cumulative sum,  $W_k(x)$ , can be used to provide both a graphical and numerical assessment of the functional form of the covariate; the cumulative sum,  $W_f(x)$ , can be used to provide an assessment of the adequacy of the link function. An omnibus test (“supremum” test) of the adequacy of the marginal regression model with respect to the relevant coordinate (e.g., a particular covariate or the fitted values) can be obtained by comparing the maximum absolute value of the observed cumulative sum to a large number of realizations (e.g., 10,000) from the null distribution (see Section 10.4 for additional details concerning the supremum test).

An appealing property of the graphical and numerical methods based on cumulative (and moving) sums of residuals is that they are valid regardless of whether the covariance among the responses has been correctly specified. As such, these graphical and numerical techniques for assessing the marginal model for the mean response are relatively robust to the working correlation assumption.

Residual analyses are also useful for detecting outlying *observations*. The detection of outlying observations is important because they can potentially have an inordinate influence on the analysis. In addition residuals can be used to identify outlying *individuals* who have unusual patterns of responses. Specifically, for each individual we can calculate a summary measure of multivariate distance between their vectors of observed and fitted responses,

$$d_i = r_i' \hat{V}_i^{-1} r_i,$$

where  $r_i$  denotes the  $n_i \times 1$  vector of residuals for the  $i^{th}$  subject,

$$V_i = A_i^{\frac{1}{2}} \text{Corr}(Y_i) A_i^{\frac{1}{2}},$$

and  $A_i^{\frac{1}{2}}$  is a diagonal matrix with the standard deviations,  $\sqrt{\phi v(\mu_{ij})}$ , along the diagonal. If the model for the mean is correctly specified, and  $V_i \approx \text{Cov}(Y_i)$ ,  $d_i$  has an approximate chi-squared distribution with degrees of freedom equal to  $n_i$ . Outlying individuals will have distances,  $d_i$ , that have relatively small associated  $p$ -values. These  $p$ -values provide a common metric for comparing  $d_i$  when the number of repeated measurements varies across subjects; however, it must be recognized that relatively small  $p$ -values (e.g.,  $p$ -values less than 0.05) are expected to occur with predictable regularity.

Finally, both influential observations and influential individuals can be detected using “deletion diagnostics” for marginal regression models fitted by generalized estimating equations. Various measures of influence, initially developed for standard linear regression (e.g., Cook’s distance and leverage), have been extended to marginal regression models. For example, Cook’s distance is a widely used statistic in standard linear regression (Cook, 1977) that measures the change in the regression parameter estimates caused by deleting each observation in turn. Cook’s distance can also be applied to marginal regression models to detect either influential individuals or influential observations. For detecting influential individuals, Cook’s distance is based on  $(\hat{\beta} - \hat{\beta}_{[i]})$ , where  $\hat{\beta}$  is the vector of parameter estimates obtained from the analysis of the data on all individuals and  $\hat{\beta}_{[i]}$  denotes the vector of parameter estimates obtained after deleting all of the observations on the  $i^{th}$  individual. For detecting influential observations, Cook’s distance is based on  $(\hat{\beta} - \hat{\beta}_{[ij]})$ , where  $\hat{\beta}_{[ij]}$  denotes the vector of parameter estimates obtained after deleting the  $j^{th}$  observation on the  $i^{th}$  individual. Exact values for these diagnostics can be obtained by deleting the observation (or the set of correlated observations on an individual) and re-fitting the marginal model to the remaining observations. However, this requires iterating the GEE fitting algorithm until convergence has been achieved. Therefore, to greatly reduce the computational burden, one-step approximations to these diagnostics are commonly used. These one-step approximations are sufficiently accurate for most practical purposes. Other measures of influence, such as leverage, have also been extended to

marginal regression models, and these statistics can be defined at the levels of individuals and observations. For example, leverage defined at the level of an individual is simply the sum of the observation leverages,  $\sum_{j=1}^{n_i} h_{ij}$ .

In summary, the GEE estimator  $\hat{\beta}$  provides a valid estimate of  $\beta$  when the model for the mean response has been correctly specified. To assess the adequacy of the marginal regression model residual analysis techniques, similar to those described in Chapter 10, can be used. In addition many of the standard regression diagnostics for identifying outliers and influential observations have been extended to marginal models fit by generalized estimating equations.

## 13.4 CASE STUDIES

Next we illustrate the main ideas presented in this and the previous chapter by considering marginal models for analyzing longitudinal data from three different studies. The first illustration employs marginal models to analyze obesity data in a sample of school-age children from the Muscatine Coronary Risk Factor (MCRF) study. The second illustration considers marginal models for analyzing count data from a study comparing two antibiotics to a placebo for the treatment of leprosy. The third illustration considers marginal models for analyzing treatment-related changes in an ordinal response measuring patients' self-assessment of their arthritis.

# Muscatine Coronary Risk Factor Study

The Muscatine Coronary Risk Factor (MCRF) study was a longitudinal survey of school-age children in Muscatine, Iowa (Woolson and Clarke, 1984; Lauer et al., 1997). The goal of the study was to examine the development and persistence of risk factors for coronary disease in children. In the MCRF study, weight and height measurements of five cohorts of children, initially aged 5–7, 7–9, 9–11, 11–13, and 13–15 years, were obtained biennially from 1977 to 1981. In total, data were collected on 4856 boys and girls. Although each child was eligible to participate in all three surveys, the data are incomplete for many children.

In this section we present longitudinal analyses of a binary response, indicating whether the child is obese. At each occasion, on the basis of a comparison of their weight to age-gender-specific norms, children were classified as obese or not obese. The goal of the analyses is to determine whether the risk of obesity increases with age and whether patterns of change in obesity are the same for boys and girls. The percentages of the children classified as obese at each of the three measurement occasions are displayed in [Table 13.1](#). These percentages were calculated from the available data in each age-gender cohort at each occasion. These descriptive statistics suggest that the rates of obesity increase from ages 6 to 12, but decline thereafter. They also suggest that the rates of obesity are higher for girls at all ages.

**Table 13.1** Percentage of children from the Muscatine Coronary Risk Factor study classified as obese in 1977, 1979, and 1981.

Gender	Age Cohort	Percentage Obese		
		1977	1979	1981
<b>Males</b>				
5–7	7.9	15.4	21.2	
7–9	18.8	20.5	23.7	
9–11	21.2	22.7	22.5	
11–13	24.3	21.8	19.4	
13–15	19.2	21.1	18.2	
<b>Females</b>				
5–7	14.0	17.2	25.1	
7–9	16.5	24.0	24.9	
9–11	25.4	26.2	22.2	
11–13	23.8	22.1	19.9	
13–15	22.9	25.8	20.9	

Initially our analysis of these data assumes that there are no cohort effects. The marginal probability of obesity is modeled as a logistic function of the covariates: linear and quadratic age, gender, and the gender-age interactions. Here age is the midpoint of the age cohort that a child belongs to (e.g., 6, 8, and 10 years at the first, second, and third occasions for the cohort of children initially aged 5–7 years). Letting  $Y_{ij} = 1$  if the  $i^{th}$  child is classified as obese at the  $j^{th}$  occasion, and  $Y_{ij} = 0$  otherwise, we assume that the marginal probability of obesity at each occasion follows the logistic model,

$$\log \left\{ \frac{\Pr(Y_{ij} = 1)}{\Pr(Y_{ij} = 0)} \right\} = \beta_1 + \beta_2 \text{Gender}_i + \beta_3 \text{Age}_{ij} + \beta_4 \text{Age}_{ij}^2 + \beta_5 \text{Gender}_i \times \text{Age}_{ij} + \beta_6 \text{Gender}_i \times \text{Age}_{ij}^2,$$

where  $\text{Age}_{ij}$  = midpoint of age cohort at the  $j^{th}$  occasion – 12 years;  $\text{Gender}_i = 1$  if the  $i^{th}$  child is female, and  $\text{Gender}_i = 0$  otherwise. This specifies the first component of a marginal model, the model for the mean response. It is assumed that the log odds of obesity changes curvilinearly with age (i.e., quadratic age trend), but the trend over time is allowed to be different for girls and boys. Next we assume that

$$\text{Var}(Y_{ij}) = \mu_{ij}(1 - \mu_{ij}).$$

This specifies the second component, the variance function and known scale parameter ( $\phi = 1, j = 1, \dots, 3$ ). Finally, we need to make assumptions about the pairwise within-subject associations among the binary responses. Because the response is binary, correlation is not the most appealing metric for association. As was mentioned in Section 12.2, with binary responses the correlations are constrained and must satisfy certain linear inequalities determined by the marginal probabilities. Instead, we specify the association in terms of pairwise log odds ratios, a more natural measure of association between pairs of binary responses.<sup>1</sup> Specifically, the within-subject association among the three repeated binary responses is assumed to have the following unstructured pairwise log odds ratio pattern,

$$\log \text{OR}(Y_{ij}, Y_{ik}) = \alpha_{jk},$$

where

$$\text{OR}(Y_j, Y_k) = \frac{\Pr(Y_j = 1, Y_k = 1)}{\Pr(Y_j = 1, Y_k = 0)} \cdot \frac{\Pr(Y_j = 0, Y_k = 0)}{\Pr(Y_j = 0, Y_k = 1)}.$$

The estimated regression coefficients and pairwise log odds ratios for the within-subject association obtained using the GEE approach are presented in [Table 13.2](#). A test of the hypothesis that changes in the log odds of obesity are the same for boys and girls,  $H_0: \beta_5 = \beta_6 = 0$ , can be constructed using a multivariate Wald statistic. This test produces a Wald statistic,  $W^2 = 0.91$ , with 2 df ( $p > 0.60$ ), and the null hypothesis cannot be rejected at the 0.05 significance level. Thus a marginal logistic regression model without the gender  $\times$  age interactions is defensible. Note that the  $\hat{\beta}$  have interpretation in terms of the pairwise log odds ratio for the responses at the  $j^{th}$  and  $k^{th}$  occasions. The pairwise log odds ratios between adjacent occasions are very similar and approximately equal to 3, indicating that the odds ratio for within-subject association is approximately 20 (or  $e^3$ ). As expected, there is strong positive association among the indicators of obesity status at the three measurement occasions.

**Table 13.2** Parameter estimates and standard errors from marginal logistic regression model for the obesity data from the Muscatine Coronary Risk Factor study.

Variable	Estimate	SE <sup>a</sup>	Z
Intercept	-1.2135	0.0506	-24.00
Gender	0.1159	0.0711	1.63
Age	0.0378	0.0133	2.85
Age <sup>2</sup>	-0.0175	0.0034	-5.19
Gender $\times$ Age	0.0075	0.0182	0.41
Gender $\times$ Age <sup>2</sup>	0.0039	0.0046	0.85
$\alpha_{12}$	3.1528	0.1280	24.63
$\alpha_{13}$	2.5975	0.1353	19.20
$\alpha_{23}$	2.9868	0.1236	24.17

<sup>a</sup>SE based on “sandwich” variance estimator.

Recall that these data on obesity are from five cohorts of children, initially aged 5–7, 7–9, 9–11, 11–13, and 13–15 years. Our analysis of the trends in the risk of obesity with age implicitly assumes that there are no cohort effects. That is, the logistic model for the probability of obesity assumes that the cross-sectional and longitudinal effects of aging are identical (see Chapter 9, Section 9.5). Following the approach used in the analysis of the FEV<sub>1</sub> data in Section 9.6, we can conduct a formal test of equality of the cross-sectional and longitudinal effects of aging by including linear and quadratic effects of both mean age (where averaging is over time) and current age minus mean age (and also their interactions with gender) in the logistic model,

$$\begin{aligned} \log \left\{ \frac{\Pr(Y_{ij} = 1)}{\Pr(Y_{ij} = 0)} \right\} = & \beta_1 + \beta_2 \text{Gender}_i + \beta_3 \overline{\text{Age}}_i + \beta_4 \overline{\text{Age}}_i^2 + \beta_5 \text{Gender}_i \times \overline{\text{Age}}_i \\ & + \beta_6 \text{Gender}_i \times \overline{\text{Age}}_i^2 + \beta_7 (\text{Age}_{ij} - \overline{\text{Age}}_i) + \beta_8 (\text{Age}_{ij}^2 - \overline{\text{Age}}_i^2) \\ & + \beta_9 \text{Gender}_i \times (\text{Age}_{ij} - \overline{\text{Age}}_i) + \beta_{10} \text{Gender}_i \times (\text{Age}_{ij}^2 - \overline{\text{Age}}_i^2), \end{aligned}$$

where  $\overline{\text{Age}}_i = \frac{1}{3} \sum_{j=1}^3 \text{Age}_{ij}$  and  $\overline{\text{Age}_i^2} = \frac{1}{3} \sum_{j=1}^3 \text{Age}_{ij}^2$ . This model distinguishes between the cross-sectional effects of aging ( $\beta_3, \beta_4, \beta_5$ , and  $\beta_6$ ) and the longitudinal effects of aging ( $\beta_7, \beta_8, \beta_9$ , and  $\beta_{10}$ ). Note that, when  $\beta_3 = \beta_7, \beta_4 = \beta_8, \beta_5 = \beta_9$ , and  $\beta_6 = \beta_{10}$ , we obtain the logistic model considered previously. A test of equality of the cross-sectional and longitudinal effects of aging,

$$H_0: (\beta_3 - \beta_7) = (\beta_4 - \beta_8) = (\beta_5 - \beta_9) = (\beta_6 - \beta_{10}) = 0,$$

produces a (multivariate) Wald statistic,  $W^2 = 2.06$ , with 4 df ( $p > 0.70$ ). This suggests that the results for aging presented in [Table 13.2](#) are not confounded by cohort effects.

Next we consider a marginal logistic regression model without the gender  $\times$  age interactions. Specifically, we consider the following model for the log odds of obesity:

$$\log \left\{ \frac{\Pr(Y_{ij} = 1)}{\Pr(Y_{ij} = 0)} \right\} = \beta_1 + \beta_2 \text{Gender}_i + \beta_3 \text{Age}_{ij} + \beta_4 \text{Age}_{ij}^2,$$

while retaining the same assumptions about the variances and pairwise log odds ratios. The estimated regression coefficients (and pairwise log odds ratios) for this model are presented in [Table 13.3](#). The estimated effect of age<sup>2</sup> is significant at the 0.05 level and these results provide evidence that the log odds of obesity increases from 6 to 12 years, levels off between age 12 to age 14, and declines between 14 to 18 years. Although the rates of obesity are significantly higher for girls at all ages, the patterns of change in the rates of obesity over time do not depend on gender. To translate these results on to a more interpretable scale, we can estimate the probability of obesity at each age for boys and girls,

**Table 13.3** Parameter estimates and standard errors from marginal logistic regression model for the obesity data from the Muscatine Coronary Risk Factor study, omitting gender  $\times$  age and gender  $\times$  age<sup>2</sup> interactions.

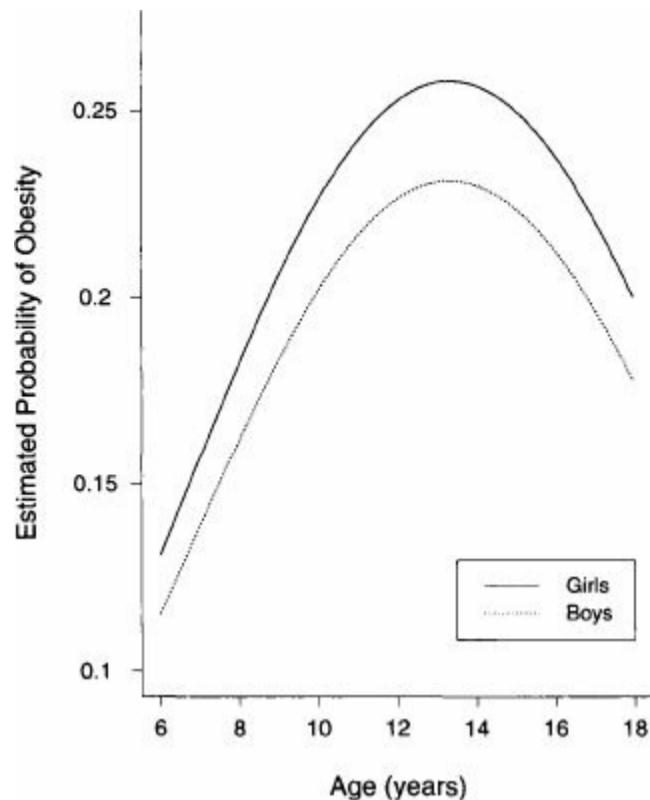
Variable	Estimate	SE <sup>a</sup>	Z
Intercept	-1.2283	0.0477	-25.75
Gender	0.1449	0.0627	2.31
Age	0.0418	0.0091	4.60
Age <sup>2</sup>	-0.0155	0.0023	-6.73
$\alpha_{12}$	3.1496	0.1280	24.61
$\alpha_{13}$	2.5931	0.1352	19.17
$\alpha_{23}$	2.9878	0.1236	24.18

<sup>a</sup>SE based on "sandwich" variance estimator.

$$\hat{\mu}_{ij} = \frac{e^{\hat{\beta}_1 + \hat{\beta}_2 \text{Gender}_i + \hat{\beta}_3 \text{Age}_{ij} + \hat{\beta}_4 \text{Age}_{ij}^2}}{1 + e^{\hat{\beta}_1 + \hat{\beta}_2 \text{Gender}_i + \hat{\beta}_3 \text{Age}_{ij} + \hat{\beta}_4 \text{Age}_{ij}^2}}.$$

For example, the estimated probability of obesity for boys at ages 6, 10, 14, and 18 is 0.12, 0.20, 0.23, and 0.18, respectively; for girls, the estimated probability of obesity at ages 6, 10, 14, and 18 is 0.13, 0.22, 0.26, and 0.20, respectively (see [Figure 13.1](#)). Note that with the logistic model, an additive effect of gender does not translate into a constant difference over time in the probability of obesity. Potential confounding of these trends by cohort effects can be examined by including linear and quadratic effects of both mean age (where averaging is over time) and current age minus mean age in the logistic model. A test of equality of the cross-sectional and longitudinal effects of aging produces a (multivariate) Wald statistic,  $W^2 = 1.74$ , with 2 df ( $p > 0.40$ ), suggesting that the results for aging presented in [Table 13.3](#) are not confounded by cohort effects.

**Fig. 13.1** Estimated probability of obesity versus age for boys and girls in the Muscatine Coronary Risk Factor study.



For illustrative purposes, we consider the same model for the log odds of obesity,

$$\log \left\{ \frac{\Pr(Y_{ij} = 1)}{\Pr(Y_{ij} = 0)} \right\} = \beta_1 + \beta_2 \text{Gender}_i + \beta_3 \text{Age}_{ij} + \beta_4 \text{Age}_{ij}^2,$$

and retain the same assumptions about the variances, but assume a Toeplitz pattern for the log odds ratios:

$$\log \text{OR}(Y_{ij}, Y_{ij+1}) = \alpha_1,$$

$$\log \text{OR}(Y_{ij}, Y_{ij+2}) = \alpha_2.$$

The estimated parameters for this model are displayed in [Table 13.4](#), and the estimates of the regression parameters (and their standard errors) are very similar to those reported in [Table 13.3](#). The estimates of the within-subject log odds ratios display a characteristic decreasing time-dependence as the time separation increases. Finally, since the Toeplitz pattern for the within-subject log odds ratios is nested within the unstructured pairwise log odds ratio model, it is possible to assess the goodness of fit of the Toeplitz model. That is, by appropriately reparameterizing the unstructured pairwise log odds ratio model,

**Table 13.4** Parameter estimates and standard errors from marginal logistic regression model for the obesity data from the Muscatine Coronary Risk Factor study, omitting gender  $\times$  age and gender  $\times$  age $^2$  interactions.

Variable	Estimate	SE <sup>a</sup>	Z
Intercept	-1.2270	0.0477	-25.72
Gender	0.1445	0.0627	2.31
Age	0.0416	0.0091	4.58
Age $^2$	-0.0156	0.0023	-6.77
$\alpha_1$	3.0684	0.0957	32.07
$\alpha_2$	2.5929	0.1353	19.17

<sup>a</sup>SE based on "sandwich" variance estimator.

$$\log \text{OR}(Y_{i1}, Y_{i2}) = \alpha_1,$$

$$\log \text{OR}(Y_{i1}, Y_{i3}) = \alpha_2,$$

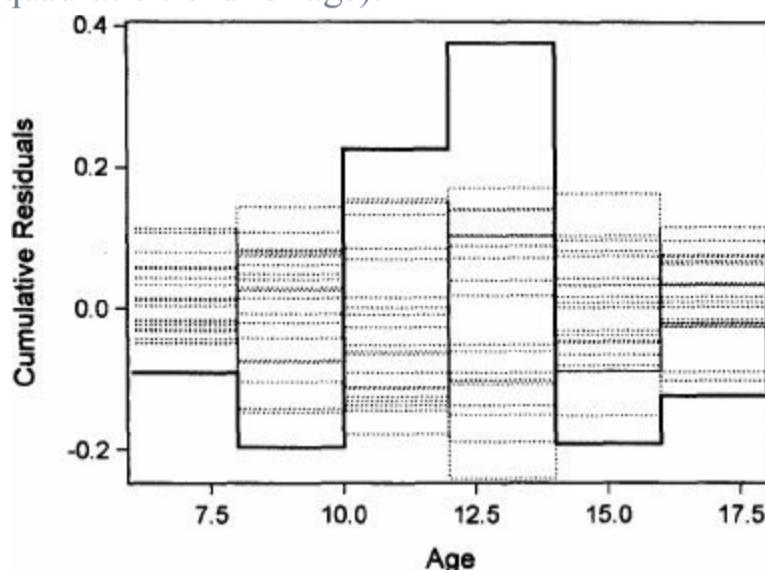
$$\log \text{OR}(Y_{i2}, Y_{i3}) = \alpha_1 + \alpha_2,$$

a 1-degree-of-freedom goodness-of-fit test (based on a Wald test) for the Toeplitz model can be constructed. The test of the null hypothesis,  $H_0: \alpha_3 = 0$ , produces a Wald  $Z = -0.99$  ( $p > 0.30$ ), indicating that the Toeplitz pattern is defensible for these data.

Finally, we consider the use of cumulative sums of residuals (see Sections 10.4 and 13.3) to assess the adequacy of the model for the probability of obesity. The results presented so far have assumed

that the log odds of obesity changes curvilinearly with age, allowing the rates of obesity to increase from 6 to 12 years, level off between age 12 to age 14, and decline between 14 to 18 years (see [Figure 13.1](#)). Specifically, the curvilinear trend in the log odds of obesity is assumed to be a quadratic function of age. [Figure 13.2](#) shows a plot of the observed cumulative sum of the residuals (solid curve), with respect to age for the quadratic trend model. On the vertical axis is the cumulative sum of residuals; the horizontal axis denotes age (in years). Superimposed on the graph are 20 simulated realizations (dotted curves) of the cumulative sum from the null distribution under the assumption that the model for the log odds of obesity is correctly specified. The realizations of the cumulative sum under the null are computer simulated from the appropriate Gaussian mean-zero process. By comparing the observed cumulative sum to the 20 different realizations under the null, it is possible to determine whether any apparent trend is systematic or due to chance fluctuations; a more formal comparison is made using the  $p$ -value for the supremum test based on 10,000 simulated realizations from the null distribution.

**Fig. 13.2** Plot of observed cumulative sum of residuals versus age (solid curve) and 20 simulated realizations (dotted curves) from the null distribution assuming a correctly specified model for the log odds of obesity (with quadratic trend for age).



From [Figure 13.2](#) it would appear that the observed cumulative sum displays a systematic pattern. In particular, the observed cumulative sum is far too large between 10 and 14 years; it also appears to be too small before 10 years and after 14 years. This suggests that the assumed functional form for age, a quadratic trend, may not be adequate. This graphical assessment of fit can be complemented by a numerical assessment. The maximum absolute value of the observed cumulative sum is 0.376. The supremum test yields a  $p$ -value of 0.0039, based on 10,000 simulated realizations of the process under the null. That is, out of 10,000 simulated realizations under the null hypothesis that a quadratic trend is adequate, only 39 had a maximum absolute value that exceeded 0.376. Thus both the graphical and numerical results suggest that the functional form for age may be inappropriate.

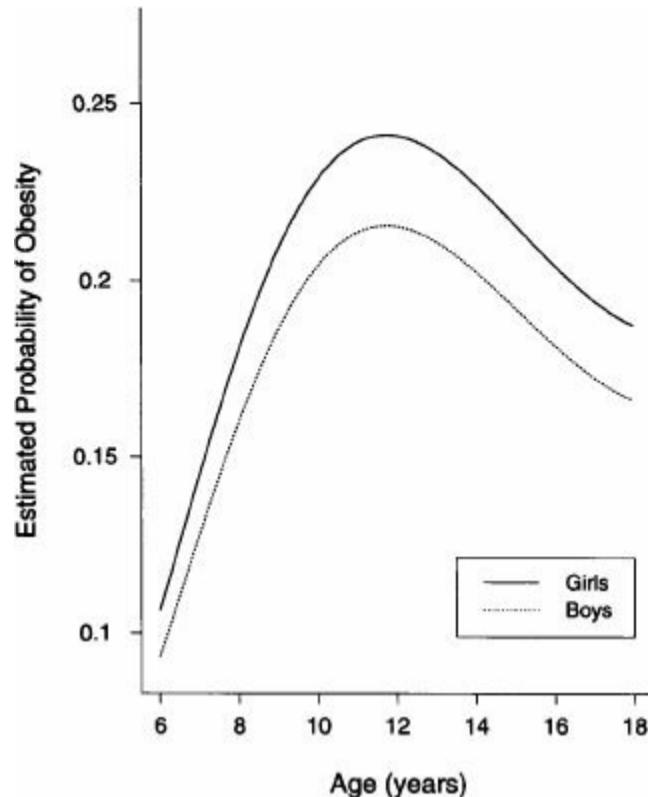
Next we consider a refinement to the model to allow for a cubic trend in the log odds as a function of age. Specifically, we consider the following model for the log odds of obesity:

$$\log \left\{ \frac{\Pr(Y_{ij} = 1)}{\Pr(Y_{ij} = 0)} \right\} = \beta_1 + \beta_2 \text{Gender}_i + \beta_3 \text{Age}_{ij} + \beta_4 \text{Age}_{ij}^2 + \beta_5 \text{Age}_{ij}^3,$$

while retaining the same assumptions about the variances and pairwise log odds ratios as before. The estimated regression coefficients (and pairwise log odds ratios) for this model are presented in [Table 13.5](#). The estimated effect of age<sup>3</sup> is significant at the 0.05 level ( $Z = 3.01$ ,  $p < 0.003$ ). These results provide evidence that the log odds of obesity departs from a quadratic trend in age. To translate these results onto a more interpretable scale, we can plot the estimated probability of obesity at each age for boys and girls (see [Figure 13.3](#)). From [Figure 13.3](#) it is clear that the rates of obesity increase sharply from 6 to 12 years but then decline, albeit at a slower rate, thereafter. Although the rates of obesity are significantly higher for girls at all ages, the overall pattern of change in the log odds of obesity does not depend on gender (a 3 df test of interaction with gender yields a chi-square statistic of 0.82,  $p > 0.80$ ). Overall, these results suggest that the rates of obesity level off earlier and decline at a somewhat less steep rate than was indicated by the quadratic trend

model (compare [Figures 13.1](#) and [13.3](#)).

**Fig. 13.3** Estimated probability of obesity versus age for boys and girls in the Muscatine Coronary Risk Factor study.



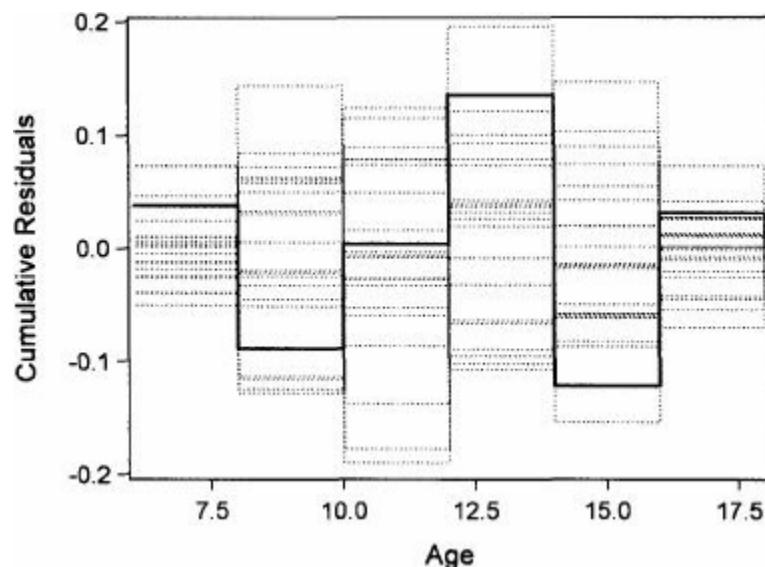
**Table 13.5** Parameter estimates and standard errors from marginal logistic regression model for the obesity data from the Muscatine Coronary Risk Factor study, assuming a cubic trend for age.

Variable	Estimate	SE <sup>a</sup>	Z
Intercept	-1.2228	0.0477	-25.65
Gender	0.1457	0.0627	2.33
Age	0.0078	0.0144	0.54
Age <sup>2</sup>	-0.0166	0.0024	-6.99
Age <sup>3</sup>	0.0018	0.0006	3.01
$\alpha_{12}$	3.1501	0.1290	24.42
$\alpha_{13}$	2.6135	0.1353	19.32
$\alpha_{23}$	2.9933	0.1231	24.31

<sup>a</sup>SE based on "sandwich" variance estimator.

We can assess the adequacy of this revised model using cumulative sums of residuals. [Figure 13.4](#) shows a plot of the observed cumulative sum of the residuals, with respect to age; superimposed on the graph are 20 realizations from the Gaussian mean-zero null distribution. This plot suggests there is no systematic trend in the observed curve. This is confirmed by a numerical assessment. The maximum absolute value of the observed cumulative sum is 0.135, with corresponding *p*-value for the supremum test equal to 0.346. With the inclusion of a cubic trend for age, both the graphical and numerical diagnostics no longer provide evidence of model misspecification.

**Fig. 13.4** Plot of observed cumulative sum of residuals versus age (solid curve) and 20 simulated realizations (dotted curves) from the null distribution assuming a correctly specified model for the log odds of obesity (with cubic trend for age).



In summary, the results indicate that the rates of obesity are significantly higher for girls at all ages. However, the overall pattern of change in the log odds of obesity does not depend on gender. For both boys and girls the rates of obesity increase sharply from 6 to 12 years but then decline, albeit at a slower rate, thereafter.

# Clinical Trial of Antibiotics for Leprosy

Next we consider count data from a placebo-controlled clinical trial of 30 patients with leprosy at the Eversley Childs Sanitarium in the Philippines (Snedecor and Cochran, 1967). Participants in the study were randomized to either of two antibiotics (denoted by treatment drug A and B) or to a placebo (denoted by treatment drug C). Prior to receiving treatment, baseline data on the number of leprosy bacilli at six sites of the body where the bacilli tend to congregate were recorded for each patient. After several months of treatment, the number of leprosy bacilli at the six sites of the body were recorded a second time. The outcome variable is the total count of the number of leprosy bacilli at the six sites.

Before proceeding with the analysis, a feature of these data should be noted. These data display substantially greater variability than that predicted by the mean under a Poisson distribution assumption. The mean number of bacilli and the variance are displayed in [Table 13.6](#). Although the sample sizes are relatively small, these descriptive statistics reveal that the variances are substantially greater than the means. As a result a Poisson assumption for the variance, with  $\text{Var}(Y_{ij}) = \mu_{ij}$ , is not appropriate for these data. Instead, we consider

**Table 13.6** Mean count of leprosy bacilli at six sites of the body (and variance) pre- and post-treatment.

Treatment Group	Baseline	Post-Treatment
Drug A (Antibiotic)	9.3 (22.7)	5.3 (21.6)
Drug B (Antibiotic)	10.0 (27.6)	6.1 (37.9)
Drug C (Placebo)	12.9 (15.7)	12.3 (51.1)

$$\text{Var}(Y_{ij}) = \phi \mu_{ij},$$

where it is assumed that  $\phi > 1$ .

In this study, the question of main scientific interest is whether treatment with antibiotics (drugs A and B) reduces the abundance of leprosy bacilli at the six sites of the body when compared to placebo (drug C). To address this question we can compare the changes, from baseline to follow-up, in the average count of leprosy bacilli in the three treatment groups. This can be expressed in the following marginal model for the expected counts of leprosy bacilli

$$\log E(Y_{ij}) = \log \mu_{ij} = \beta_1 + \beta_2 \text{Time}_{ij} + \beta_3 \text{Time}_{ij} \times \text{Trt}_{1i} + \beta_4 \text{Time}_{ij} \times \text{Trt}_{2i},$$

where  $Y_{ij}$  is the count of the number of leprosy bacilli for the  $i^{th}$  patient in the  $j^{th}$  period of observation ( $j = 1, 2$ ). The variables  $\text{Trt}_1$  and  $\text{Trt}_2$  are indicator variables for drugs A and B respectively, with  $\text{Trt}_1 = 1$  if a patient was randomized to drug A and  $\text{Trt}_1 = 0$  otherwise, and  $\text{Trt}_2 = 1$  if a patient was randomized to drug B and  $\text{Trt}_2 = 0$  otherwise. The binary variable,  $\text{Time}$ , denotes the baseline and post-treatment follow-up periods, with  $\text{Time} = 0$  for the baseline period (period 1) and  $\text{Time} = 1$  for the post-treatment follow-up period (period 2). Because patients were randomized to one of the three treatments, the model does not include main effects of treatment (since the mean count of the number of leprosy bacilli at baseline can be assumed to be equal in the three treatment groups). To complete the specification of the model, we must make assumptions about the variances of the counts and the within-subject association among the repeated counts. Because of the discernible *overdispersion* in these data (relative to Poisson variability), we assume that the variance of  $Y_{ij}$  is given by

$$\text{Var}(Y_{ij}) = \phi \mu_{ij},$$

where  $\phi$  can be thought of as an overdispersion factor. Finally, the within-subject association is accounted for by assuming a common correlation,

$$\text{Corr}(Y_{i1}, Y_{i2}) = \alpha.$$

In this marginal model for the expected number of leprosy bacilli, all of the covariates are

dichotomous and the log-linear regression parameters can be given interpretations in terms of (log) rate ratios. In [Table 13.7](#) we summarize the interpretation of  $\beta$  in terms of the log expected counts in the three groups at baseline and during post-treatment follow-up. So, for example, the expected count of leprosy bacilli at the six sites of the body at baseline in the placebo group (drug C) is  $e^{\beta_1}$ , while the expected count during the follow-up period is  $e^{\beta_1 + \beta_2}$ . Thus  $e^{\beta_2}$  is the rate ratio of leprosy bacilli, comparing the follow-up period to baseline, in the placebo group (drug C). Similarly  $e^{\beta_2 + \beta_3}$  is the rate ratio of leprosy bacilli, comparing the follow-up period to baseline, in the group randomized to drug A. Finally,  $e^{\beta_2 + \beta_4}$  is the rate ratio of leprosy bacilli, comparing the follow-up period to baseline, in the group randomized to drug B.

**Table 13.7** Parameters of the marginal log-linear regression model for the leprosy bacilli data.

Treatment Group	Period	$\log(\mu_{ij})$
Drug A (Antibiotic)	Baseline	$\beta_1$
	Follow-up	$\beta_1 + \beta_2 + \beta_3$
Drug B (Antibiotic)	Baseline	$\beta_1$
	Follow-up	$\beta_1 + \beta_2 + \beta_4$
Drug C (Placebo)	Baseline	$\beta_1$
	Follow-up	$\beta_1 + \beta_2$

As a result a direct comparison of the three treatment groups in terms of changes in the expected rates of leprosy bacilli is expressible in terms of  $\beta_3$  and  $\beta_4$ . That is,  $\beta_3$  and  $\beta_4$  represent the difference between the changes in the log expected rates, comparing drug A and B to the placebo (drug C). For example, a value of  $\beta_3 < 0$  indicates a greater reduction in the rate of bacilli from baseline in the group randomized to drug A (when compared to the placebo group).

The estimated regression coefficients, obtained using the GEE approach, are displayed in [Table 13.8](#) (with standard errors based on the “sandwich” estimator). A test of the null hypothesis,  $H_0: \beta_3 = \beta_4 = 0$ , produces a (multivariate) Wald statistic,  $W^2 = 6.99$ , with 2 degrees of freedom ( $p < 0.05$ ). This indicates that treatment with antibiotics significantly reduces the abundance of leprosy bacilli at the six sites of the body. A test of the null hypothesis that the two antibiotics are equally effective,  $H_0: \beta_3 = \beta_4$ , produces a Wald statistic,  $W^2 = 0.08$ , with 1 degree of freedom ( $p > 0.7$ ). Thus we cannot reject the null hypothesis that the two antibiotics are equally effective in reducing the number of leprosy bacilli. To obtain a common estimate of the log rate ratio, comparing both antibiotics (drugs A and B) to placebo, we can fit the reduced model

**Table 13.8** Parameter estimates and standard errors (based on sandwich variance estimator) from marginal log-linear regression model for the leprosy bacilli data.

Variable	Estimate	SE	Z
Intercept	2.3734	0.0801	29.62
Time	-0.0138	0.1573	-0.09
Time × Trt <sub>1</sub>	-0.5406	0.2186	-2.47
Time × Trt <sub>2</sub>	-0.4791	0.2279	-2.10

Note: Estimated scale or dispersion parameter:  $\hat{\phi} = 3.45$ . Estimated working correlation:  $\hat{\alpha} = 0.797$ .

$$\log E(Y_{ij}) = \log \mu_{ij} = \beta_1 + \beta_2 \text{Time}_{ij} + \beta_3 \text{Time}_{ij} \times \text{Trt}_i,$$

where the variable Trt is an indicator variable for antibiotics, with Trt = 1 if a patient was randomized to either drug A or B and Trt = 0 otherwise. We retain the same assumptions about the variance and correlation as before.

The estimated regression coefficients are displayed in [Table 13.9](#) (with standard errors based on the “sandwich” estimator). The common estimate of the log rate ratio, comparing post-treatment rates of bacilli in the antibiotics group (drugs A and B) to placebo, is -0.5141. Thus the rate ratio is 0.60 (or  $e^{-0.5141}$ ), with 95% confidence interval, 0.41 to 0.88, indicating that treatment with antibiotics

significantly reduces the average number of bacilli when compared to placebo. In the placebo group, there is a non-significant reduction in the average number of bacilli of approximately 1 % (or  $[1 - e^{-0.0108}] \times 100\%$ ), while in the antibiotics group there is a significant reduction of approximately 40% (or  $[1 - e^{-0.0108-0.5141}] \times 100\%$ ).

**Table 13.9** Parameter estimates and standard errors (based on sandwich variance estimator) from marginal log-linear regression model for the leprosy bacilli data, with common estimate of the effect of antibiotics.

Variable	Estimate	SE	Z
Intercept	2.3734	0.0801	29.62
Time	-0.0108	0.1572	-0.07
Time × Trt	-0.5141	0.1966	-2.62

*Note:* Estimated scale or dispersion parameter:  $\hat{\phi} = 3.41$ .

Finally, the estimated pairwise correlation is relatively large (approximately 0.8), suggesting that there may be substantial heterogeneity among patients in their disease severity. Of note, the estimated scale parameter is approximately 3.4, revealing overdispersion relative to that predicted by Poisson variability. Recall from Section 13.2 that in addition to correcting for misspecification of the within-subject association, the “sandwich” estimator of the standard errors also corrects for any misspecification of the variance, including overdispersion. Therefore standard errors based on the “sandwich” estimator are automatically adjusted for potential overdispersion; it is not necessary to include an additional parameter,  $\phi$ , to account for overdispersion. It is instructive to compare the standard errors in [Table 13.9](#) with the corresponding model-based standard errors when  $\phi$  is fixed at 1 (Poisson variance assumption) and when  $\phi$  is estimated from the data (allowing for overdispersion by a constant factor  $\phi > 1$ ). Results for the former are presented in [Table 13.10](#); results for the latter are presented in [Table 13.11](#). First, note that regardless of whether  $\phi$  is fixed at 1 or estimated as an additional parameter, the estimates of the regression parameters,  $\hat{\beta}$ , are the same. Second, the model-based standard errors in [Table 13.10](#) are discernibly smaller than those in [Table 13.9](#), reflecting the fact that they are based on the assumption of Poisson variability. Third, the model-based standard errors in [Table 13.11](#) are larger than those reported in [Table 13.10](#) by a constant factor of 1.845 (or  $\sqrt{3.41}$ ). Although the standard errors in [Tables 13.11](#) and [13.9](#) have both made adjustments for overdispersion, the former corrections are based on a constant multiple of the standard errors obtained under a Poisson variance assumption whereas the latter corrections allow for more general departures from Poisson variability. For example, standard errors based on the “sandwich” estimator make corrections for overdispersion that may be differential by treatment group. Finally, we note that overdispersion relative to Poisson variation can also be accounted for by assuming that the counts have negative binomial variance. Assuming negative binomial variance for the counts of leprosy bacilli (and a log link function) yields estimates of the regression parameters and standard errors that are qualitatively very similar to those in [Table 13.11](#). Recall that the negative binomial variance allows for overdispersion by assuming the variance increases as a *quadratic* function of the mean (see Section 11.5). However, as with the use of a constant scale factor (which implicitly assumes the variance increases as a *linear* function of the mean), this correction for overdispersion does not allow for more general departures from Poisson variability. In contrast, corrections based on the “sandwich” variance estimator allow for *any* departures from Poisson variability when the conditions required for its use are met (see Section 13.2).

**Table 13.10** Parameter estimates and model-based standard errors from marginal log-linear regression model for the leprosy bacilli data, with common estimate of the effect of antibiotics.

Variable	Estimate	SE	Z
Intercept	2.3734	0.0557	42.59
Time	-0.0108	0.0619	-0.17
Time × Trt	-0.5141	0.0832	-6.18

Note: Fixed scale or dispersion parameter:  $\phi = 1$ .

**Table 13.11** Parameter estimates and model-based standard errors from marginal log-linear regression model for the leprosy bacilli data, with common estimate of the effect of antibiotics.

Variable	Estimate	SE	Z
Intercept	2.3734	0.1028	23.08
Time	-0.0108	0.1142	-0.09
Time × Trt	-0.5141	0.1536	-3.35

Note: Estimated scale or dispersion parameter:  $\hat{\phi} = 3.41$ .

# Arthritis Clinical Trial

The final example is from a longitudinal clinical trial comparing auranofin therapy (3 mg of oral gold, twice daily) and placebo for the treatment of rheumatoid arthritis (Bombardier et al., 1986). In this six-month, randomized, double-blind trial, 303 patients with classic or definite rheumatoid arthritis were randomized to one of the two treatment groups and followed over time. The outcome variable of interest is a global impression scale (Arthritis Categorical Scale) measured at baseline (month 0), month 2, month 4, and month 6. This is a self-assessment of a patient's current arthritis, measured on a five-level ordinal scale: (1) very good, (2) good, (3) fair, (4) poor, and (5) very poor. Baseline data on this outcome variable are available for 303 of the patients who participated in this trial; follow-up data at 6 months are available for 294 patients.

The goal of the analysis is to assess changes in the odds of a more favorable response over the duration of the study, and to determine whether treatment with auranofin has an influence on these changes. Letting  $Y_{ij}$  denote the ordinal response for the  $i^{th}$  subject at the  $j^{th}$  occasion, we assume that the log odds of a more favorable response at each occasion follows the proportional odds model

$$\log \left\{ \frac{\Pr(Y_{ij} \leq k)}{\Pr(Y_{ij} > k)} \right\} = \alpha_k + \beta_1 \text{Trt}_i + \beta_2 \sqrt{\text{Month}_{ij}} \\ + \beta_3 \text{Trt}_i \times \sqrt{\text{Month}_{ij}},$$

where  $\text{Month}_{ij}$  is the timing of the measurement, in months, for the  $i^{th}$  subject at the  $j^{th}$  occasion,  $\text{Trt}_i = 1$  if the  $i^{th}$  subject is randomized to auranofin, and  $\text{Trt}_i = 0$  if randomized to placebo. This specifies the first component of a marginal model, the model for the mean response. Specifically, the model assumes that the log odds of a favorable response changes linearly with (square-root transformed) time, but the slopes over time are allowed to differ between the two treatment groups. Preliminary analyses indicated that changes in the log odds of a more favorable response were approximately linear in square-root transformed time; a 2 df Wald test of non-linearity (quadratic trend) yielded  $W^2 = 0.18$ ,  $p > 0.90$ . The second component of the marginal model, the variance of the multinomial response, is completely determined by the model for the mean response. For the third component, we initially make a “working independence” assumption for the within-subject association among the repeated ordinal responses and rely on the empirical variance estimator for making valid inferences.

The GEE estimates in [Table 13.12](#) indicate that the trajectories for the log odds over time are significantly different for patients treated with placebo versus patients treated with auranofin ( $Z = 2.66$ ,  $p < 0.01$ ). Specifically, patients treated with auranofin have a significantly greater increase in the odds of a more favorable response over the duration of the study. Relative to baseline the odds of a more favorable response at month 6 has increased by a factor of 1.84 (or  $e^{0.2481 \times \sqrt{6}}$ ) for the placebo group but by a factor of 3.27 (or  $e^{(0.2481+0.2354) \times \sqrt{6}}$ ) for the auranofin group. At the completion of the trial, patients treated with auranofin are approximately twice (or  $e^{0.2354 \times \sqrt{6}} = 1.78$ ) as likely to have a more favorable response when compared to patients treated with placebo. Not surprisingly, due to the randomization,  $\hat{\beta}_1 \approx 0$  indicating that the two groups have a similar log odds of a favorable response at baseline (month 0).

**Table 13.12** GEE estimates and standard errors (empirical) from the proportional odds model for the arthritis clinical trial data.

Variable	Estimate	SE	Z
$\alpha_1$	-3.1902	0.1994	-16.00
$\alpha_2$	-1.2042	0.1523	-7.91
$\alpha_3$	0.5736	0.1464	3.92
$\alpha_4$	2.4770	0.1995	12.42
Trt	0.0714	0.1975	0.36
$\sqrt{\text{Month}}$	0.2481	0.0613	4.05
Trt $\times$ $\sqrt{\text{Month}}$	0.2354	0.0883	2.66

Finally, for illustrative purposes we re-fit the model with an unstructured pattern for the within-subject association among the repeated ordinal responses. In general, specification of a “working covariance” for ordinal responses (other than “working independence”) is very challenging because it requires the specification and estimation of a large number of parameters. With  $n$  repeated measures of a  $K$ -level ordinal response, there are  $(K - 1)^2 \times n \times (n - 1)/2$  pairwise parameters. In our example with four repeated measures of a five-level ordinal response, an unstructured pattern has 96 pairwise association parameters that require estimation. Fitting a model with an unstructured pattern yielded estimates and standard errors very similar to those reported in [Table 13.12](#). For the effect of main interest, the  $\text{Trt} \times \sqrt{\text{Month}}$  interaction, the analysis with unstructured pattern for the within-subject association yields  $\hat{\beta}_3 = 0.2377$  (model-based SE = 0.0869, empirical SE = 0.0877). This estimate of  $\hat{\beta}_3$  is similar to the estimate reported in [Table 13.12](#) ( $\hat{\beta}_{43} = 0.2354$ ) under a “working independence” assumption for the within-subject association. In addition the empirical SE for  $\hat{\beta}_3$  reported in [Table 13.12](#) (empirical SE = 0.0883) is very similar to both the model-based and empirical standard errors obtained from the analysis with unstructured pattern for the within-subject association. The fact that the empirical standard errors are so similar under these two “working” assumptions for the within-subject association suggests that there has been negligible loss of efficiency from basing the analysis on a “working independence” assumption; however, we caution that this cannot be expected in general.

## 13.5 MARGINAL MODELS AND TIME-VARYING COVARIATES

In this section<sup>†</sup> we return to an implicit assumption in marginal models that was highlighted in the previous chapter at the end of Section 12.2. Marginal models assume that the conditional mean of the  $j^{th}$  response, given  $X_{i1}, \dots, X_{in_i}$ , depends only on  $X_{ij}$ . As we will see, this assumption has some important ramifications for time-varying covariates in marginal models. Recall that the vector of covariates at the  $j^{th}$  occasion,  $X_{ij}$ , includes two main types of covariates: covariates whose values do not change throughout the duration of the study and covariates whose values change over time. The former are referred to as time-invariant or between-subject covariates (e.g., gender and fixed experimental treatments), while the latter are referred to as time-varying or within-subject covariates (e.g., time since baseline, current smoking status, and environmental exposures). In this section we consider estimation and aspects of interpretation of time-varying covariates in marginal models. We want to emphasize at the outset that the issues discussed here in the context of marginal models apply equally to the linear models for longitudinal data in Part II of the book.

When considering time-varying covariates, we can distinguish covariates that vary systematically over time but are fixed by design of the study and covariates that vary randomly over time. An example of a time-varying covariate that is fixed by design is a treatment group indicator in a crossover trial. Another example, and one that is commonly encountered in a longitudinal study, is time since baseline (when the measurement occasions are fixed by the study design). Covariates that vary randomly over time are often referred to as *stochastic*, that is, values of the covariate at any occasion cannot be precisely predicted since they are governed by a random mechanism. An example of a time-varying covariate that is stochastic is current blood glucose level. In an observational study of diabetics, participants' blood sugar levels can vary randomly over the duration of the study. Additional examples include current smoking status or cumulative pack-years, blood pressure, cholesterol level, fat intake, and exposure to environmental pollutants. As we will later see, when a covariate is both time-varying and stochastic, new issues arise concerning the interpretation and estimation of regression parameters in marginal models for longitudinal data.

Marginal models for the mean response described in this and earlier chapters can be specified as

$$(13.4) \quad g(\mu_i) = g\{E(Y_i|X_i)\} = X_i\beta,$$

for some known link function  $g(\cdot)$ . This use of vector and matrix notation implies that the model for the mean at each occasion is given by

$$g(\mu_{ij}) = g\{E(Y_{ij}|X_i)\} = X'_{ij}\beta, \quad (j = 1, \dots, n_i).$$

However, what is often overlooked is the implicit assumption that the conditional mean of the  $j^{th}$  response, given  $X_{i1}, \dots, X_{in_i}$ , depends only on  $X_{ij}$

$$(13.5) \quad E(Y_{ij}|X_i) = E(Y_{ij}|X_{i1}, \dots, X_{in_i}) = E(Y_{ij}|X_{ij}).$$

With time-invariant covariates, this assumption necessarily holds since  $X_{ij} = X_{ik}$  for all occasions  $k \neq j$ . Also, with time-varying covariates that are fixed by design of the study (e.g., treatment group indicator in a crossover trial), the assumption also holds since values of the covariates at any occasion are determined a priori by study design and in a manner completely unrelated to the longitudinal response. However, when a covariate is time-varying and stochastic (13.5) may not necessarily hold. For example, the assumption will be violated when the current value of  $Y_{ij}$ , given  $X_{ij}$ , predicts the subsequent value of  $X_{ij+1}$ . In that case

$$E(Y_{ij}|X_{ij}, X_{ij+1}) \neq E(Y_{ij}|X_{ij}),$$

and  $X_{ij+1}$  is said to confound the relationship between  $Y_{ij}$  and  $X_{ij}$ . In general, when (13.5) does not hold, then preceding and/or subsequent values of the time-varying covariate confound the relationship between  $Y_{ij}$  and  $X_{ij}$ ; this can lead to biased estimates of  $\beta$  in the marginal model given by (13.4).

To fix ideas, consider a longitudinal study designed to examine the effects of physical exercise on reducing blood glucose levels in patients with type 2 diabetes mellitus. We let  $X_{ij}$  denote the cumulative amount of physical activity at the  $j^{th}$  occasion and  $Y_{ij}$  denote a measure of blood glucose. The goal of the study is to determine the relationship between  $Y_{ij}$  and  $X_{ij}$ . Next suppose that subjects with elevated blood glucose levels at the  $j^{th}$  occasion subsequently increase their level of physical activity, while subjects with the same cumulative amount of physical activity at the  $j^{th}$  occasion, but with normal blood glucose levels, continue to maintain their usual level of physical activity. Then the assumption that

$$E(Y_{ij}|X_{i1}, \dots, X_{in_i}) = E(Y_{ij}|X_{ij})$$

does not hold and the relationship between  $Y_{ij}$  and  $X_{ij}$  is confounded by  $X_{ij+1}$ . In particular, the strength of the relationship between  $Y_{ij}$  and  $X_{ij}$  will be underestimated if subjects with elevated blood glucose levels subsequently increase their amount of physical activity.

In general, when a covariate is time-varying and stochastic, much greater care is needed in modeling its relationship to the response variable. It is important to assess the assumption made in (13.5), namely that the conditional mean of  $Y_{ij}$ , given the entire time-varying covariate profile  $X_{i1}, \dots, X_{in_i}$ , depends only on the covariate value at the  $j^{th}$  occasion,  $X_{ij}$ . Note, however, that  $X_{ij}$  can be defined in terms of functions of the explanatory variables measured at or preceding the  $j^{th}$  occasion (e.g., cumulative exposure at the  $j^{th}$  occasion). When (13.5) is violated the relationship between the mean of  $Y_{ij}$  and  $X_{ij}$ , expressed in terms of  $\beta$ , will be confounded by preceding and/or subsequent values of the covariate and misleading inferences about  $\beta$  can result.

Finally, even when (13.5) holds, there can be problems with the *interpretation* of regression parameters relating the mean response to stochastically time-varying covariates. In particular, the regression parameters  $\beta$  in (13.4) may not have the implied causal interpretation. For example, the model given by (13.4) may correctly specify the relationship between mean blood glucose level and physical activity at the last measurement occasion, since at the last occasion  $X_{in_i}$ , the cumulative amount of physical activity, is a function of the entire time-varying covariate profile. However, even though (13.5) holds, the regression parameters  $\beta$  may not have the implied causal interpretation without making additional assumptions. To see why, let us consider a simplified version of the example discussed earlier.

Suppose that a group of diabetics are measured at two occasions. Let  $Y_{i1}$  and  $Y_{i2}$  denote the blood glucose levels at baseline and follow-up, and  $X_{i1}$  and  $X_{i2}$  denote measures of physical activity at the two occasions. Suppose that it is of interest to determine the association between the cumulative amount of physical activity,  $X_i^* = X_{i1} + X_{i2}$ , and blood glucose level at the completion of the study,  $Y_{i2}$ . The following model is assumed:

$$E(Y_{i2}|X_i^*) = \beta_1 + \beta_2 X_i^*,$$

where, for ease of exposition, an identity link function is assumed. In this model  $\beta_2$  appears to have interpretation as the effect of a unit increase in the cumulative amount of physical activity on the mean blood glucose level at follow-up, since

$$E(Y_{i2}|X_i^* = x + 1) - E(Y_{i2}|X_i^* = x) = \beta_2.$$

However, because  $X_{ij}$  is time-varying and stochastic, this interpretation of  $\beta_2$  rests on the validity of *either* of the following two assumptions: (1)  $Y_{i2}$  is not predicted by  $Y_{i1}$ , given  $X_{i1}$  and  $X_{i2}$ , or (2)  $X_{i2}$  is not predicted by  $Y_{i1}$ , given  $X_{i1}$ . In particular, if neither of these assumptions holds,  $Y_{i1}$  “confounds” the relationship between  $Y_{i2}$  and  $X_i^*$  and  $\beta_2$  does not have the desired causal interpretation. We loosely use the term “confounding” to emphasize that  $Y_{i1}$  obscures or distorts the association of real scientific interest between  $Y_{i2}$  and  $X_i^*$ . Strictly speaking,  $Y_{i1}$  can be considered both a “confounder” and an “intermediate variable” on the causal path between  $Y_{i2}$  and  $X_i^*$ . When  $Y_{i1}$  is both a confounder and an intermediate variable, standard methods of adjustment for confounding no longer apply (since  $Y_{i1}$  is predicted by  $X_{i1}$ , and so should not be adjusted for, but also predicts  $X_{i2}$ , and so should be

adjusted for in the analysis of the association between  $Y_{i2}$  and  $X_i^*$ ). Instead, advanced statistical methods for causal inference (e.g., marginal structural models and structural nested models; see references at end of chapter) are required when neither assumption 1 nor 2 holds. However, a discussion of statistical methods for causal inference is beyond the scope of this chapter; some references to the statistical literature on this topic appear at the end of the chapter.

Let us consider these two assumptions in context. In a longitudinal study, it is unlikely that assumption 1 would ever hold, since the repeated responses are usually positively correlated (given the covariates,  $X_{i1}$  and  $X_{i2}$ ). Therefore the causal interpretation of  $\beta_2$  usually rests on the validity of assumption 2. For example, assumption 2 would be violated if subjects with elevated blood glucose levels at baseline subsequently increase their level of physical activity, while subjects with the same amount of physical activity at baseline, but with normal blood glucose levels, continue to maintain their usual level of physical activity. When assumption 2 holds, the covariate is said to be *external* with respect to the response variable and  $\beta$  has the desired causal interpretation.

In summary, when a covariate is both time-varying and stochastic, we must consider the relationship between the response at any occasion, say  $Y_{ij}$ , and the subsequent value of the covariate,  $X_{ij+1}$ . A time-varying covariate is said to be *external* when the current and preceding values of the response at the  $j^{th}$  occasion ( $Y_{i1}, \dots, Y_{ij}$ ), given the current and preceding values of the time-varying covariate ( $X_{i1}, \dots, X_{ij}$ ), do not predict the subsequent value of  $X_{ij+1}$ . More formally, a time-varying covariate is *external* (or sometimes referred to as *exogenous*) when

$$(13.6) \quad f(X_{ij+1}|X_{i1}, \dots, X_{ij}, Y_{i1}, \dots, Y_{ij}) = f(X_{ij+1}|X_{i1}, \dots, X_{ij});$$

otherwise, the covariate is said to be *internal* (or *endogenous*). This generalizes assumption 2. Note that when a covariate is external,

$$E(Y_{ij}|X_i) = E(Y_{ij}|X_{i1}, \dots, X_{in_i}) = E(Y_{ij}|X_{i1}, \dots, X_{ij}),$$

which is a weaker assumption than (13.5). An example of an external covariate is air pollution in studies of children's lung function growth. The outdoor levels of air pollutants (e.g., ozone, fine suspended particulate matter, and sulfur dioxide) are time-varying and stochastic, but conditional on past values, future values are not predicted by the lung function responses of the study participants and (13.6) holds. Note, however, that children's personal exposure to air pollution would not be considered an external covariate if children with poor lung function growth subsequently altered their daily behavior (e.g., spending less time outdoors) to avoid exposure to high levels of air pollution. In principle, it is possible to examine the assumption that a time-varying covariate is *external* by considering regression models for the dependence of  $X_{ij}$  on  $Y_{i1}, \dots, Y_{ij-1}$  (or some known function(s) of  $Y_{i1}, \dots, Y_{ij-1}$ ) and  $X_{i1}, \dots, X_{ij-1}$  (or some known function(s) of  $X_{i1}, \dots, X_{ij-1}$ ). The absence of any relationships between  $X_{ij}$  and  $Y_{i1}, \dots, Y_{ij-1}$ , given the preceding covariate profile,  $X_{i1}, \dots, X_{ij-1}$ , provides support for the validity of the assumption that the covariate process is *external*.

In conclusion, when covariates are time-varying and stochastic the regression parameters do not necessarily have the implied causal interpretation even when (13.5) holds. The regression parameters can be given a causal interpretation only when it can be further assumed that the time-varying covariates are external with respect to the response variable (i.e., when (13.6) holds).

# 13.6 COMPUTING: GENERALIZED ESTIMATING EQUATIONS USING PROC GENMOD IN SAS

To fit marginal models using the generalized estimating equations approach, we can use an enhanced option for repeated measures data in the PROC GENMOD procedure in SAS. Although PROC GENMOD is primarily a procedure for fitting generalized linear models to a single response, the use of a REPEATED statement in PROC GENMOD allows for the fitting of marginal models to correlated responses using the GEE approach. In Chapter 15, Section 15.6, we describe how an alternative procedure in SAS, PROC GLIMMIX, can also fit marginal models using the GEE approach.

For example, to fit a marginal logistic regression model to longitudinal data from two groups, with the within-subject associations specified in terms of log odds ratios, we can use the illustrative SAS commands given in [Table 13.13](#). Similarly, to fit a marginal log-linear regression model to longitudinal data from two groups, with the within-subject associations specified in terms of correlations, we can use the illustrative SAS commands given in [Table 13.14](#). To fit a marginal proportional odds model to longitudinal ordinal data, under a “working independence” assumption for the within-subject association, we can use the illustrative SAS commands given in [Table 13.15](#). To assess the adequacy of the functional form for the time trend in the marginal model for the mean based on cumulative sums of residuals, we can use the illustrative SAS commands given in [Table 13.16](#). Next we describe the most salient parts of the command syntax required for fitting marginal models to longitudinal data using the GEE approach within PROC GENMOD in SAS.

**Table 13.13** Illustrative commands for a marginal logistic regression, with within-subject associations specified in terms of log odds ratios, using PROC GENMOD in SAS.

---

```
PROC GENMOD DESCENDING;
  CLASS id group time;
  MODEL y=group time group*time/DIST=BINOMIAL LINK=LOGIT;
  REPEATED SUBJECT=id/WITHINSUBJECT=time LOGOR=FULLCLUST;
```

---

**Table 13.14** Illustrative commands for a marginal log-linear regression, with within-subject associations specified in terms of correlations, using PROC GENMOD in SAS.

---

```
PROC GENMOD;
  CLASS id group time;
  MODEL y=group time group*time/DIST=POISSON LINK=LOG;
  REPEATED SUBJECT=id/WITHINSUBJECT=time TYPE=UN;
```

---

**Table 13.15** Illustrative commands for a marginal proportional odds regression, with a “working independence” assumption for the within-subject associations, using PROC GENMOD in SAS.

---

```
PROC GENMOD;
  CLASS id group;
  MODEL y=group time group*time/DIST=MULT LINK=CUMLOGIT;
  REPEATED SUBJECT=id/TYPE=IND;
```

---

**Table 13.16** Illustrative commands for requesting model assessment of the functional form for the time trend based on cumulative sums of residuals, with a “working independence” assumption for the within-subject associations, using PROC GENMOD in SAS.

---

```
PROC GENMOD;
  CLASS id group;
  MODEL y=group time group*time / DIST=BINOMIAL LINK=LOGIT;
  REPEATED SUBJECT=id / TYPE=IND;
```

ASSESS VAR(time) / RESAMPLES = 10000 SEED=7435865;

---

PROC GENMOD <options>;

This statement calls the procedure GENMOD in SAS. It can also include an option for specifying the level of the response variable that is modeled. By default, the lower response level is modeled. For a binary response, coded (0,1), it is the probability that  $Y = 0$  that is modeled. For an ordinal response, coded (1, 2, ..., K), the response categories are ordered from lowest to highest and the probabilities of the lower response levels are modeled. Use of the DESCENDING option reverses the default ordering of the response levels, resulting in the highest response level(s) being modeled (i.e., the probability that  $Y = 1$  for binary data that are coded as 0 and 1).

CLASS variables;

The CLASS statement is used to define all variables that are to be regarded as categorical or factors. By default, this statement will create indicator variables for each factor using a reference group coding, with the last level (where “last” here refers to the level with the largest alphanumeric value) regarded as the reference group. Different sort orders for the CLASS variables can be requested by the ORDER=<option> on the PROC GENMOD statement.

MODEL response = <effects> / <options>;

MODEL events/trials = <effects> / <options>;

The MODEL statement specifies the response variable and the covariate effects. The second form of the MODEL statement, with the events/trials syntax, allows the response to be in the form of a ratio of two variables (e.g., counts of the number of successes and the number of trials) and is used for binomial response data. The linear predictor can include both discrete (defined in the CLASS statement) and quantitative (excluded from the CLASS statement) covariates. By default, PROC GENMOD includes a column of 1's for the intercept in the model.

The option that ordinarily is used to specify the distribution of a single univariate response has a somewhat different role when fitting a marginal model using the GEE approach. The option DIST=*keyword* does not specify a distribution for the vector of correlated responses; instead, it specifies the default canonical link function and variance function that happen to be associated with particular exponential family distributions. For example, the option DIST=POISSON does not specify that the response vector (or even its separate components) has a Poisson distribution; instead, it specifies that the mean of the response vector is related to the covariates via a log link function (the canonical link for the Poisson distribution), and the mean and variance of the responses are related by  $\text{Var}(Y) = E(Y) = \mu$  (i.e., the variance function is  $v(\mu) = \mu$ ).

Note that PROC GENMOD also provides a wide choice of options for the inclusion of a dispersion parameter,  $\phi$ . However, the scale parameter  $\phi$  is assumed to be time-invariant. This restriction on the scale parameter is a limitation of the implementation of the GEE approach that makes it unappealing for analyzing longitudinal data when the response variable is continuous and the variance of the repeated measurements is not constant over the duration of the study (e.g., the variance of baseline measurements is often discernibly smaller than post-baseline measurements).

The LINK=*keyword* specifies the choice of built-in link function relating the mean response to the linear predictor. If the LINK=*keyword* is omitted, the default link function is the canonical link function associated with the particular exponential family distribution specified on DIST=*keyword*.

A final option often required when modeling count data is an offset. The OFFSET=*variable* specifies a variable to be used as an offset. For example, in modeling count data, the rate is often of more direct interest and the denominator for the counts or “population at risk” can be included as an offset. Note that this variable cannot be a CLASS variable, and it should not be included as one of the covariates listed on the MODEL statement.

REPEATED SUBJECT=subject-effect / <options>;

The REPEATED statement distinguishes the fitting of a generalized linear model for a single

univariate responses via maximum likelihood from the fitting of a marginal model to a vector of correlated responses using the GEE approach. The REPEATED statement is used to specify the assumed structure of the within-subject association among the repeated measurements.

In particular, the REPEATED statement defines a variable that determines the clustering of observations within an individual. The latter is achieved by including a subject identifier, that distinguishes clusters of correlated responses, on the SUBJECT=*subject-effect*; this is not optional, a *subject-effect* must be included with the REPEATED statement and this variable must be listed in the CLASS statement. By including a subject identifier, pairs of observations with the same value of that variable are regarded as correlated (by virtue of arising from the same subject) while pairs of observations with distinct values are regarded as independent.

A useful option on the REPEATED statement is the WITHINSUBJECT=*within-subject effect*. With this option a variable denoting the “repeated effect” can be included and this identifies the order of the repeated measurements within subjects. In the context of longitudinal data, the “repeated effect” identifies the measurement occasions. While it is not always necessary to include this variable, failure to do so may have unforeseen consequences when there are vectors of repeated measures of different length and/or when the vector of responses are not in the same order for all subjects. To avoid any potential problems, this variable should be included on the REPEATED statement, whenever possible, to ensure that the within-subject association is estimated appropriately.

While the REPEATED statement in PROC GENMOD has a similar function to the REPEATED statement in PROC MIXED, the order in which the *subject-effect* and the *within subject-effect* appear in the REPEATED statement are reversed (for reasons perhaps best known only to the developers at SAS Institute). By default, PROC GENMOD produces a table of regression parameter estimates, standard errors, and Z statistics. The standard errors and Z statistics are based on the empirical or “sandwich” estimator of  $\text{Cov}(\hat{\beta})$  described in Section 13.2. Use of the REPEATED statement with the MODELSE option produces the corresponding table based on the “model-based” estimator of  $\text{Cov}(\hat{\beta})$ .

Finally, two additional options are used for specifying assumptions about the structure of the working correlation matrix or the log odds ratios (for binary responses only) among the repeated measurements. The TYPE=*correlation-structure* specifies the working correlation structure. PROC GENMOD provides a number of build-in correlation structures, including unstructured (UN),  $m$ -dependent (MDEP( $m$ ), where  $m$  is the order of dependence), first-order autoregressive (AR), and exchangeable (analogous to “compound symmetry”) or equicorrelated (EXCH/CS). For ordinal responses with a multinomial distribution, PROC GENMOD currently only supports a “working independence” assumption (IND); indeed, one of the main challenges with extending the GEE approach to ordinal responses has been that the “working covariance” for ordinal responses (other than “working independence”), in general, requires the specification and estimation of a large number of nuisance parameters. For binary responses only, the structure of the within-subject association among the responses can be specified in terms of log odds ratios using the LOGOR=*log odds ratio structure* option. For example, in [Table 13.13](#), the LOGOR=FULLCLUST option estimates separate log odds ratios for all pairs of responses; this is analogous to an “unstructured” odds ratio pattern. PROC GENMOD also allows a very flexible regression structure for the log odds ratios. Note that either the TYPE=*correlation-structure* or the LOGOR=*log odds ratio structure* option should be specified, but not both. By default, a working independence structure is assumed.

**ASSESS** VAR=*effect* | LINK/<options>;

The ASSESS statement computes and plots statistics based on aggregates of residuals. Three types of aggregates are available: cumulative sums of residuals, moving sums of residuals, and lowess smoothed residuals; the default is cumulative sums. To create an analysis, either VAR=*effect* or LINK must be specified. VAR=*effect* requests that the functional form of a covariate be assessed by performing the analysis with respect to the variable identified by the effect. The effect must be specified in the MODEL statement and must be a continuous variables.

LINK requests the assessment of the link function by performing the analysis with respect to the linear predictor.

The WINDOW and LOESS options in the ASSESS statement requests model assessment based on moving sums of residuals and lowess smoothed residuals respectively.

An important option in the ASSESS statement is the RESAMPLES<=number> option; this specifies the number of paths used for computing the *p*-value for the supremum test (the default is 1,000 simulated paths). Another useful option is the SEED=number option; this specifies a seed for the normal random number generator used in creating simulated realizations of aggregates of residuals for plots and estimating *p*-values. Specifying a seed allows you to reproduce identical graphs and *p*-values from a later run of the procedure.

Of note, the initial output produced by PROC GENMOD is the standard output from a generalized linear model assuming that all observations are independent. The resulting estimates of the regression coefficients are used as initial values for the generalized estimating equations algorithm. However, the reader is cautioned that this initial output should be ignored. In particular, the reported value of the log-likelihood and various likelihood-based goodness of fit statistics should not be considered part of the GEE output.

## **13.7 FURTHER READING**

Burton et al. (1998) provide an accessible introduction to generalized estimating equations. A more comprehensive description of generalized estimating equations can be found in Chapter 6 of the textbook by Myers et al. (2001); also see Chapter 11 of the textbook by Agresti (2002).

# Bibliographic Notes

The early foundations for statistical methods for the analysis of repeated categorical responses can be traced to a general approach developed by Grizzle, Starmer, and Koch (1969); this approach became known as the GSK method. Koch et al. (1977) applied the GSK method to the analysis of repeated measurements. However, the application of the GSK method was limited to categorical covariates. The GEE approach overcame many of the limitations of the GSK method.

The theoretical foundation for the generalized estimating equations approach can be found in Godambe (1960) and Durbin (1960); also see Huber (1967, 1981) and White (1982). Liang and Zeger (1986) and Zeger and Liang (1986), in companion papers, proposed a class of generalized estimating equations for repeated measures and longitudinal data; see Liang and Zeger (1995) for a historical perspective on generalized estimating equations. Connections between the GEE approach and likelihood-based methods were made by Zhao et al. (1992), Fitzmaurice and Laird (1993), and Fitzmaurice et al. (1993).

The “sandwich” variance estimator was suggested by Cox (1961) and derived in Huber (1967), White (1982), Gourieroux et al. (1984), and Royall (1986); see Hinkley and Wang (1991) and Kauermann and Carroll (2001) for a discussion of properties of the “sandwich” variance estimator. For finite samples, simulation studies have shown that Wald tests using the sandwich estimator tend to be liberal, that is, have nominal  $p$ -values that are too small (see Lin and Wei, 1989; Emrich and Piedmonte, 1992; Gunsolley et al., 1995; Fay et al., 1998; Mancl and DeRouen, 2001; Fay and Graubard, 2001).

The implicit assumption in marginal regression models with time-varying covariates given by (13.5) is discussed in Fitzmaurice et al. (1993), Pepe and Anderson (1994), Robins et al. (1999), Pan et al. (2000), and Pan and Connett (2002); also see Schildcrout and Heagerty (2005) for a discussion of the bias–variance trade-off when the assumption does not hold.

A general discussion of methods for estimating the causal effect of time-varying covariates in marginal models for longitudinal data can be found in Robins et al. (1999) and the references therein; also see Chapter 23 of Fitzmaurice et al. (2009). Chapter 12 of Diggle et al. (2002) presents a useful summary of the key ideas.

Regression diagnostics for marginal models fit by generalized estimating equations were developed by Preisser and Qaqish (1996). They provide computational formulae for one-step approximations to deletion diagnostics for the influence of a single observation and for the influence of the set of correlated observations on a single individual (or cluster).

## Problems

**13.1** In a clinical trial of patients with respiratory illness, 111 patients from two different clinics were randomized to receive either placebo or an active treatment. Patients were examined at baseline and at four visits during treatment. At each examination, respiratory status (categorized as 1 = good, 0 = poor) was determined. These data are from Koch et al. (1990), and are reported in Davis (1991) and Stokes et al. (1995). The main objective of the analyses is to understand the joint effects of treatment and time on the probability that respiratory status is classified as good. It is also of interest to determine whether the effect of treatment is the same for patients from the two clinics.

The raw data are stored in an external file: `respir.dat`

Each row of the data set contains the following eight variables:

ID Clinic Treatment Y<sub>0</sub> Y<sub>1</sub> Y<sub>2</sub> Y<sub>3</sub> Y<sub>4</sub>

*Note:* The respiratory status response variable  $Y_j$  is coded 1 = good, and 0 = poor, at the  $j^{th}$  occasion.

The categorical (character) variable Treatment is coded A = Active drug, P = Placebo. The categorical variable Clinic is coded 1 = clinic 1, 2 = clinic 2.

**13.1.1** Ignoring the clinic variable, consider a model for the log odds that respiratory status is classified as good, including the main effects of treatment and time (where time is regarded as a

categorical variable with five levels), and their interaction.

Use generalized estimating equations (GEE), assuming separate pairwise log odds ratios (or separate pairwise correlations, if available software does not permit the within-subject association to be parameterized in terms of log odds ratios) among the five binary responses. Construct a test of the null hypothesis of no effect of treatment on *changes* in the log odds that respiratory status is classified as good based on the empirical standard errors.

**13.1.2** What conclusions do you draw about the effect of treatment on changes in the log odds? Provide results that support your conclusions.

**13.1.3** Patients in this trial were drawn from two separate clinics. Repeat the analysis for Problem 13.1.1, allowing the effects of treatment (and, possibly, time) to depend on clinic.

(a) Is the effect of treatment the same in the two clinics? Present results to support your conclusion.

(b) Find a parsimonious model that describes the effects of clinic, treatment, and time, on the log odds that respiratory status is classified as good. For the model selected, give a clear interpretation of the estimated regression parameters for the final model selected.

**13.1.4** For the final model selected in Problem 13.1.3, construct a table of the estimated probabilities that respiratory status is classified as good as a function of both time and treatment group (and, possibly, clinic). What do you conclude from this table?

**13.2** In a clinical trial of patients suffering from epileptic seizures (Thall and Vail, 1990), patients were randomized to receive either a placebo or the drug progabide, in addition to standard therapy. A baseline count of the number of epileptic seizures in an 8-week period prior to randomization was obtained. In addition counts of the number of epileptic seizures in each of four successive 2-week (post-baseline) treatment periods were obtained. The goal of the analysis is to make a comparison between the two treatment groups in terms of changes in the rates of epileptic seizures throughout the duration of the study.

The raw data are stored in an external file: `epilepsy.dat`

Each row of the data set contains the following eight variables:

ID Y<sub>1</sub> Y<sub>2</sub> Y<sub>3</sub> Y<sub>4</sub> Treatment Y<sub>0</sub> Age

*Note:* The response variable  $Y_0$  is a baseline count of the number of epileptic seizures in an 8-week interval. The response variables  $Y_j$  are counts of the number of epileptic seizures in the four successive 2-week (post-baseline) treatment intervals, for  $j = 1, \dots, 4$ . The categorical variable Treatment is coded 1 = Progabide, 0 = Placebo. The variable Age is the age of each patient (in years) at baseline.

**13.2.1** Consider a model for the log seizure rate that includes the main effects of treatment and time (where time is regarded as a categorical variable with five levels), and their interaction. Use generalized estimating equations (GEE), assuming separate pairwise correlations among the five responses. Construct a test of the null hypothesis of no effect of treatment on changes in the log seizure rate based on the empirical standard errors.

**13.2.2** What conclusions do you draw about the effect of treatment on *changes* in the log seizure rate?

**13.2.3** Construct a new variable, Ptime, where:

Ptime = 0 if baseline, and Ptime = 1 if post-baseline (any of the four successive 2-week intervals).

Repeat the analysis for Problem 13.2.1 using Ptime (instead of time as a categorical variable with five levels). Construct a test of the null hypothesis of no effect of treatment on changes in the log seizure rate based on the empirical standard errors.

**13.2.4** From the results of the analysis for Problem 13.2.3, what conclusions do you draw about the effect of treatment on changes in the log seizure rate?

**13.2.5** Patient 49 (ID = 49) is a potential outlier. This patient reported 151 seizures during the 8-week baseline interval and 302 (102+65+72+63) seizures during the four successive 2-week

intervals. Repeat all of the analyses in Problems 13.2.1 to 13.2.4, excluding all of the repeated count data from patient 49. When the data from patient 49 are excluded, what conclusions do you draw about the effect of treatment on changes in the log seizure rate?

**13.3** In a clinical trial of patients with insomnia (Francom, Chuang-Stein, and Landis, 1989), patients were randomized to receive either a hypnotic drug or placebo. An ordinal response, denoting patients' reported time (in minutes) to fall asleep after going to bed was recorded at baseline and after two weeks of treatment. The four-level ordinal response is coded 1: < 20 minutes, 2: 20–30 minutes, 3: 30–60 minutes, 4: > 60 minutes; these data are from Chapter 11 ([Table 11.4](#)) of Agresti (2002). The main objective of the analyses is to assess changes in the odds of a more favorable response (shorter reported time to fall asleep) over the duration of the study, and to determine whether treatment with the hypnotic drug has an influence of these changes.

The raw data are stored in an external file: `insomnia.dat`

Each row of the data set contains the following four variables:

ID Trt Time Y

*Note:* The response variable  $Y$  is a four-level ordinal response denoting patients' reported time (in minutes) to fall asleep after going to bed (1: < 20 minutes, 2: 20–30 minutes, 3: 30–60 minutes, 4: > 60 minutes). The variable Time denotes baseline (Time = 0) and 2-week follow-up (Time = 1). The categorical treatment variable, Trt, is coded 1 = Hypnotic drug, 0 = Placebo. The variable Age is the age of each patient (in years) at baseline.

**13.3.1** Consider a proportional odds model for the cumulative log odds of response that includes the main effects of treatment and time, and their interaction,

$$\log \left\{ \frac{\Pr(Y_{ij} \leq k)}{\Pr(Y_{ij} > k)} \right\} = \alpha_k + \beta_1 \text{Trt}_i + \beta_2 \text{Time}_{ij} \\ + \beta_3 \text{Trt}_i \times \text{Time}_{ij}.$$

Fit the model using generalized estimating equations (GEE), with a “working independence” assumption for the within-subject association.

**13.3.2** What is the interpretation of the estimate of  $\beta_2$ ?

**13.3.3** What is the interpretation of the estimate of  $\beta_3$ ?

**13.3.4** Construct a test of the null hypothesis of no effect of treatment on changes in the cumulative log odds of response based on the empirical standard errors. What conclusions do you draw about the effect of treatment?

**13.3.5** Based on the results from Problem 13.3.1, estimate the odds ratio of a more favorable response at week 2 relative to baseline for patients receiving placebo.

**13.3.6** Based on the results from Problem 13.3.1, estimate the odds ratio of a more favorable response at week 2 relative to baseline for patients receiving the hypnotic drug.

**13.3.7** Based on the results from Problem 13.3.1, estimate the probability that patients receiving the hypnotic drug report falling asleep in less than 20 minutes (i.e., the probability of response level 1) at week 2.

**13.3.8** Based on the results from Problem 13.3.1, estimate the probability that patients receiving the hypnotic drug report falling asleep in 30–60 minutes (i.e., the probability of response level 3) at week 2.

*Hint:*  $\Pr(Y_{ij} = k) = \Pr(Y_{ij} \leq k) - \Pr(Y_{ij} \leq k - 1)$ .

<sup>1</sup> Although the (log) odds ratio is a preferable metric for within-subject association among pairs of binary responses, implementations of GEE with odds ratios for within-subject association are currently incorporated in only a few statistical software packages (e.g., the `LOGOR` option in PROC GENMOD in SAS, the `ordgee` function in the `geepack` package in R, and the `alr` package in R and S-Plus). However, because statistical software is constantly evolving, we anticipate that implementations of GEE with odds ratios for within-subject association will soon become available within most of the major statistical packages.

<sup>†</sup> This section examines in detail the implicit assumption in marginal models concerning time-varying covariates. This assumption has important implications for estimation and interpretation of time-varying covariate effects. Although the content of this section is somewhat challenging, and can be omitted on first reading without loss of continuity, we strongly encourage the reader to return to this section.

# *Chapter 14*

## ***Generalized Linear Mixed Effects Models***

### **14.1 INTRODUCTION**

In the previous chapter we described marginal models for longitudinal data. Marginal models can be considered an extension of generalized linear models that *directly* incorporate the within-subject association among the repeated measurements. To estimate the regression parameters in a marginal model, we made some assumptions about the marginal distribution of the response at each occasion (e.g., assumptions about the mean, and its dependence on the covariates, and the variance of each  $Y_{ij}$ ). We also made assumptions about the pairwise within-subject associations among the responses, thereby linking repeated observations of the same subject. A notable feature of marginal models is that the mean response and the covariance are modeled separately. This separation ensures that the interpretation of the regression coefficients in a marginal model does not rely on the assumed model for the covariance among the responses. In specifying the marginal means, variances, and pairwise associations, we did not fully specify the joint distribution of the vector of responses. However, these assumptions were sufficient for estimating and constructing confidence intervals for the regression parameters using the GEE approach.

An alternative approach for accounting for the within-subject association is via the introduction of random effects. In Chapter 8 we saw how the incorporation of random effects at the individual level induces correlation among the repeated measures at the population level. In this chapter we describe how generalized linear models can be extended to longitudinal data by allowing a subset of the regression coefficients to vary randomly from one individual to another. These models are known as *generalized linear mixed effects models*, and they extend in a natural way the conceptual approach represented by the linear mixed effects models discussed in Chapter 8. However, we must caution the reader at the outset that the introduction of random effects in generalized linear models produces a greater degree of conceptual and analytic complexity relative to marginal models or to random effects in linear models. Although both classes of models account for the within-subject association among repeated measurements, the manner in which they do so has important implications for the interpretation of the regression parameters. In Chapter 16 we highlight the major distinctions between the regression coefficients in marginal and generalized linear mixed models and consider various aspects of interpretation of the regression effects in these two classes of models for longitudinal data.

## 14.2 INCORPORATING RANDOM EFFECTS IN GENERALIZED LINEAR MODELS

The basic premise underlying the generalized linear mixed effects model for longitudinal data is the assumption of heterogeneity across individuals in the study population in a subset of the regression coefficients from a generalized linear model. That is, a subset of the regression coefficients (e.g., intercepts in a logistic regression model) are assumed to vary across individuals according to some distribution. The random effects can be thought of as reflecting natural heterogeneity due to many unmeasured factors. For mathematical and computational convenience, we ordinarily assume that the random effects have a multivariate normal distribution. Then, conditional on the random effects, we assume that the responses for any particular individual are independent observations from a distribution belonging to the exponential family (e.g., the Bernoulli distribution if  $Y_{ij}$  is binary, the Poisson distribution if  $Y_{ij}$  is a count, or the multinomial distribution if  $Y_{ij}$  is ordinal). The latter assumption is completely analogous to the “conditional independence” assumption ( $R_i = \sigma^2 I_{n_i}$ ) made in the linear mixed effects model described in Chapter 8. In fact the linear mixed effects model is simply a special case of the generalized linear mixed effects model where the conditional mean, given the random effects, is related to the covariates via an identity link function and the conditional distribution of the responses is assumed to be normal. Because the linear mixed effects model is a special case, it is useful for pedagogical purposes to consider its formulation within the framework and terminology of generalized linear models. By doing so, the extensions to other types of response variables will become more apparent.

# Linear Mixed Effects Models

In this section we consider the linear mixed effects model as a generalized linear model, albeit one with both fixed and random effects. Recall that the standard generalized linear model formulation requires a three-part specification: (1) a distributional assumption, (2) a systematic component, and (3) a link function. In the linear mixed effects model it is assumed that the *conditional* distribution of each  $Y_{ij}$ , given a vector of random effects  $b_i$ , has a normal distribution, with  $\text{Var}(Y_{ij}|b_i) = \sigma^2$  (i.e.,  $\phi = \sigma^2$  and  $v(\mu_{ij}) = 1$ ). In addition, given the random effects  $b_i$ , it is assumed that the  $Y_{ij}$  are independent of one another (i.e., given  $b_i$ ,  $Y_{ij}$  and  $Y_{ik}$  are assumed to be independent of each other). This completes the distributional assumptions on the  $Y_{ij}$ . Next the conditional mean of  $Y_{ij}$  is assumed to depend on both fixed and random effects via the following extended definition of the linear predictor:

$$\eta_{ij} = X'_{ij}\beta + Z'_{ij}b_i,$$

where the definition of the linear predictor has been extended to incorporate both population (or fixed) and subject-specific (or random) effects. In addition the random effects,  $b_i$ , are assumed to have a multivariate normal distribution. This specifies the systematic component. Finally, an identity link function relates the conditional mean of  $Y_{ij}$  to the linear predictor,

$$E(Y_{ij}|b_i) = \eta_{ij} = X'_{ij}\beta + Z'_{ij}b_i.$$

That is, for the identity link,  $\eta_{ij} - g\{E(Y_{ij}|b_i)\} = E(Y_{ij}|b_i)$  and hence,

$$E(Y_{ij}|b_i) = X'_{ij}\beta + Z'_{ij}b_i.$$

In the linear mixed effects model, the response for the  $i^{th}$  subject at the  $j^{th}$  occasion is assumed to differ from the population mean,  $X'_{ij}\beta$ , by a subject-specific effect,  $Z'_{ij}b_i$ , and a within-subject measurement error,  $\epsilon_{ij}$ . The within-subject measurement errors are independently normally distributed, with zero mean and variance  $\sigma^2$ .

When collected in a vector,  $\epsilon_i \sim N(0, R_i)$ , where  $R_i = \sigma^2 I_{n_i}$  (the “conditional independence” assumption). Recall that  $R_i = \text{Cov}(\epsilon_i)$  describes the covariance among observations when we focus on the mean response profile of any *individual*; that is, it is the covariance of the  $i^{th}$  subject’s deviations from his/her mean response profile,  $X_i\beta + Z_i b_i$ . Also the  $b_i$  are assumed to vary independently from one individual to another, with  $b_i \sim N(0, G)$ .

When expressed in vector and matrix notation, the linear mixed effects model is

$$Y_i = X_i\beta + Z_i b_i + \epsilon_i,$$

where the vector of regression parameters  $\beta$  (the fixed effects) is assumed to be the same for all individuals and the vector of subject-specific regression coefficients  $b_i$  (the random effects) describes how the  $i^{th}$  individual’s mean response profile deviates from the overall population trend. A distinctive feature of the linear mixed effects model is that it yields simple expressions for both the conditional mean response (for any individual),

$$E(Y_i|b_i) = X_i\beta + Z_i b_i,$$

and the marginal mean response (for the population), averaged over all individuals,

$$E(Y_i) = X_i\beta.$$

Thus the regression coefficients  $\beta$  have population-averaged interpretations in terms of how the mean response changes over time and how these changes relate to covariates.

Finally, the conditional covariance of the responses, given the random effects  $b_i$ , is assumed to be a diagonal matrix with

$$\text{Cov}(Y_i|b_i) = \text{Cov}(\epsilon_i) = R_i = \sigma^2 I_{n_i}.$$

On the other hand, the marginal covariance of the responses (the covariance among deviations of the  $i^{th}$  individual’s responses from the population mean,  $X_i\beta$ ),

$$\text{Cov}(Y_i) = \text{Cov}(Z_i b_i) + \text{Cov}(\epsilon_i)$$

$$= Z_i G Z_i' + R_i$$

$$= Z_i G Z_i' + \sigma^2 I_{n_i},$$

is certainly not diagonal. Thus the introduction of random effects,  $b_i$ , in the linear mixed effects model induces correlation (marginally) among the  $Y_i$ . This consequence of introducing random effects extends more generally and, in a very natural way, to any generalized linear model with random effects. That is, the correlations among the repeated observations on an individual can be thought of as arising from sharing a set of underlying random effects.

# Generalized Linear Mixed Effects Models

Next we consider how the ideas underlying the linear mixed effects model can be extended to generalized linear models. Once again, we can formulate the generalized linear mixed model using a three-part specification:

1. We assume that the conditional distribution of each  $Y_{ij}$ , given a  $q \times 1$  vector of random effects  $b_i$ , belongs to the exponential family of distributions and that  $\text{Var}(Y_{ij}|b_j) = v\{E(Y_{ij}|b_i)\} \phi$ , where  $v(\cdot)$  is a known variance function, a function of the conditional mean,  $E(Y_{ij}|b_i)$ . In addition, given the random effects  $b_i$ , it is assumed that the  $Y_{ij}$  are independent of one another; this is the “conditional independence” assumption.
2. The conditional mean of  $Y_{ij}$  is assumed to depend on fixed and random effects via the following linear predictor:

$$\eta_{ij} = X'_{ij}\beta + Z'_{ij}b_i,$$

with

$$g\{E(Y_{ij}|b_i)\} = \eta_{ij} = X'_{ij}\beta + Z'_{ij}b_i$$

for some known link function,  $g(\cdot)$ .

3. The random effects are assumed to have some probability distribution. In principle, any multivariate distribution can be assumed for the  $b_i$ , in practice, it is common to assume that the  $b_i$  have a multivariate normal distribution, with zero mean and  $q \times q$  covariance matrix,  $G$ . In addition the random effects,  $b_i$ , are assumed to be independent of the covariates,  $X_i$ .

These three components completely specify a broad class of generalized linear mixed models. Note that in Chapters 12 and 13 we extended generalized linear models by making assumptions about the mean and covariance of  $Y_i$ ; in particular, we did not make full distributional assumptions about  $Y_i$ . In contrast, the three components of a generalized linear mixed model given above completely specify the joint distribution of  $Y_i$ . To fix ideas, consider the following four illustrative examples of generalized linear mixed effects models using this three-component specification.

# Example 1: Generalized Linear Mixed Model for a Continuous Response

Suppose that  $Y_{ij}$  is a continuous response and that it is of interest to relate changes in the mean response over time to the covariates. An example of a linear mixed effects model for  $Y_{ij}$  is given by the following three-part specification:

1. Conditional on a vector of random effects  $b_i$ , the  $Y_{ij}$  are independent and assumed to have a normal distribution, with  $\text{Var}(Y_{ij}|b_i) = \sigma^2$  (i.e.,  $\phi = \sigma^2$  and the variance does not depend on the conditional mean).
2. The conditional mean of  $Y_{ij}$  depends on fixed and random effects via the following linear predictor:

$$\eta_{ij} = X'_{ij}\beta + Z'_{ij}b_i,$$

where  $X'_{ij} = Z'_{ij} - (1 \ t_{ij})$ , with

$$\begin{aligned} E(Y_{ij}|b_i) &= \eta_{ij} = X'_{ij}\beta + Z'_{ij}b_i \\ &= \beta_1 + \beta_2 t_{ij} + b_{1i} + b_{2i} t_{ij} \\ &= (\beta_1 + b_{1i}) + (\beta_2 + b_{2i})t_{ij}. \end{aligned}$$

That is, the conditional mean of  $Y_{ij}$  is related to the linear predictor by an identity link function,  $\eta_{ij} = g\{E(Y_{ij}|b_i)\} = E(Y_{ij}|b_i)$ .

3. The random effects are assumed to have a bivariate normal distribution, with zero mean and  $2 \times 2$  covariance matrix  $G$ .

This illustration of a generalized linear mixed effects model is simply a random intercepts and slopes model and is a special case of the linear mixed effects models considered in Chapter 8. However, when it is viewed as a generalized linear mixed effects model, a much broader class of models for continuous responses can, in principle, be entertained. For example, the mean can be related to the linear predictor by a link function other than the identity. Thus, if the effects of covariates are thought to act multiplicatively on the mean response, a log link function might be more appropriate. Alternatively, the variance can be allowed to depend on any known function of the mean response.

## Example 2: Generalized Linear Mixed Model for Counts

Suppose that  $Y_{ij}$  is a count. An example of a generalized linear mixed model for  $Y_{ij}$  is given by the following three-part specification:

1. Conditional on a vector of random effects  $b_i$ , the  $Y_{ij}$  are independent and have a Poisson distribution, with  $\text{Var}(Y_{ij}|b_i) = E(Y_{ij}|b_i)$ , (i.e.,  $\phi = 1$ ).
2. The conditional mean of  $Y_{ij}$  depends on fixed and random effects via the following linear predictor:

$$\eta_{ij} = X'_{ij}\beta + Z'_{ij}b_i,$$

where  $X'_{ij} = Z'_{ij} = (1, t_{ij})$ , with

$$\log \{E(Y_{ij}|b_i)\} = \eta_{ij} = X'_{ij}\beta + Z'_{ij}b_i.$$

That is, the conditional mean of  $Y_{ij}$  is related to the linear predictor by a log link function; this is an example of a log-linear mixed effects model.

3. The random effects are assumed to have a bivariate normal distribution, with zero mean and  $2 \times 2$  covariance matrix  $G$ .

In this example the model is a log-linear regression model with randomly varying intercepts and slopes. This model posits that there is natural heterogeneity among individuals in both their baseline level and changes in the expected counts over time. Note that in this example the model assumes Poisson variation for the counts, conditional on the random effects. In Section 14.4 we discuss ways to relax this assumption and allow for overdispersion relative to Poisson variability.

## Example 3: Generalized Linear Mixed Model for a Binary Response

Suppose that  $Y_{ij}$  is a binary response, taking values of 0 or 1. A logistic mixed effects model for  $Y_{ij}$  is given by the following three-part specification:

1. Conditional on a single random effect  $b_i$ , the  $Y_{ij}$  are independent and have a Bernoulli distribution, with  $\text{Var}(Y_{ij}|b_i) = E(Y_{ij}|b_i)\{1 - E(Y_{ij}|b_i)\}$ , (i.e.,  $\phi = 1$ ).
2. The conditional mean of  $Y_{ij}$  depends on fixed and random effects via the following linear predictor:

$$\eta_{ij} = X'_{ij}\beta + Z'_{ij}b_i = X'_{ij}\beta + b_i,$$

where  $Z_{ij} = 1$  for all  $i = 1, \dots, N$ , and  $j = 1, \dots, n_i$ , with

$$\log \left\{ \frac{\Pr(Y_{ij} = 1|b_i)}{\Pr(Y_{ij} = 0|b_i)} \right\} = \eta_{ij} = X'_{ij}\beta + b_i.$$

That is, the conditional mean of  $Y_{ij}$  is related to the linear predictor by a logit link function.

3. The single random effect  $b_i$  is assumed to have a univariate normal distribution, with zero mean and variance  $g_{11}$ .

In this example the model is a simple logistic regression model with randomly varying intercepts. This model can be considered a discrete data analogue of the “compound symmetry” model discussed in Chapters 7 and 8. The model posits that there is natural heterogeneity in individuals’ propensity to respond positively that persists throughout all binary responses obtained on any individual.

## Example 4: Generalized Linear Mixed Model for an Ordinal Response

Last, suppose that  $Y_{ij}$  is an ordinal response with  $K$  categories ( $1, \dots, K$ ). A logistic mixed effects model for the *cumulative response probabilities* is given by the following three-part specification:

1. Conditional on a vector of random effects  $b_i$ , the  $Y_{ij}$  are independent and have a multinomial distribution (with multinomial covariance determined by the conditional means or conditional response probabilities).
2. The  $k^{\text{th}}$  cumulative response probability for  $Y_{ij}$  depends on fixed and random effects via the following linear predictor:

$$\eta_{ij} = X'_{ij}\beta + Z'_{ij}b_i,$$

where  $X'_{ij} = Z'_{ij} = (1, t_{ij})$ , with

$$\log\{E(Y_{ij}|b_i)\} = \eta_{ij} = X'_{ij}\beta + Z'_{ij}b_i.$$

That is, the conditional cumulative response probabilities are related to the linear predictor by a logit link function.

3. The random effects are assumed to have a bivariate normal distribution, with zero mean and  $2 \times 2$  covariance matrix  $G$ .

In this example the model is a proportional odds regression model with randomly varying intercepts and slopes. This model posits that there is natural heterogeneity among individuals in both their baseline level and changes in the cumulative response probabilities over time.

Although in the first three examples we have chosen canonical link functions to relate the conditional mean of  $Y_{ij}$  to  $\eta_{ij}$ , in principle, any suitable link function can be chosen. The four examples of generalized linear mixed effects models considered so far are purely illustrative. They demonstrate how the choices of the three components might differ according to the type of response variable. However, these four examples should not be considered prescriptions for constructing generalized linear mixed effects models.

## 14.3 INTERPRETATION OF REGRESSION PARAMETERS

Although the introduction of random effects can simply be thought of as a means of accounting for the correlation among longitudinal responses, it has important implications for the interpretation of the regression coefficients in generalized linear mixed models. The regression parameters,  $\beta$ , have somewhat different interpretations than the regression parameters in the marginal models considered in Chapters 12 and 13. In generalized linear mixed models the regression coefficients have subject-specific interpretations. That is, they represent the influence of covariates on a *specific* subject's mean response. In particular, the regression coefficients are interpreted in terms of the effects of *within-subject* changes in covariates on changes in an individual's transformed mean response, while holding the remaining covariates constant. This interpretation for  $\beta$  can be better appreciated by considering the following simple example of a logistic regression model with randomly varying intercepts:

$$\log \left\{ \frac{\Pr(Y_{ij} = 1|b_i)}{\Pr(Y_{ij} = 0|b_i)} \right\} = X'_{ij}\beta + b_i,$$

where  $b_i$  is assumed to have a univariate normal distribution, with zero mean and variance  $g_{11}$ . The interpretation of a component of  $\beta$ , say  $\beta_k$ , is in terms of changes in any given *individual's* log odds of response for a unit *within-subject* change in the corresponding covariate, say  $X_{ijk}$ . That is, when  $X_{ijk}$  takes on some value  $x$ , the log odds of a positive response is

$$\begin{aligned} \log \left\{ \frac{\Pr(Y_{ij} = 1|b_i, X_{ij1}, \dots, X_{ijk} = x, \dots, X_{ijp})}{\Pr(Y_{ij} = 0|b_i, X_{ij1}, \dots, X_{ijk} = x, \dots, X_{ijp})} \right\} \\ = b_i + \beta_1 X_{ij1} + \dots + \beta_k x + \dots + \beta_p X_{ijp}. \end{aligned}$$

Similarly, when  $X_{ijk}$  now takes on some value  $x + 1$ , but all other covariate values are held fixed, the log odds of a positive response is

$$\begin{aligned} \log \left\{ \frac{\Pr(Y_{ij} = 1|b_i, X_{ij1}, \dots, X_{ijk} = x + 1, \dots, X_{ijp})}{\Pr(Y_{ij} = 0|b_i, X_{ij1}, \dots, X_{ijk} = x + 1, \dots, X_{ijp})} \right\} \\ = b_i + \beta_1 X_{ij1} + \dots + \beta_k(x + 1) + \dots + \beta_p X_{ijp}. \end{aligned}$$

Thus, for any individual, the log odds of a positive response for a unit increase in  $X_{ijk}$  is simply  $\beta_k$ . (obtained by subtracting the former log odds from the latter). That is,  $\beta_k$  measures the change in the log odds of response per unit increase in  $X_{ijk}$ , for any given individual having some unobservable underlying propensity to respond positively,  $b_i$ . Also note that this subject-specific interpretation of  $\beta_k$  is far more natural for a covariate that varies within an individual (i.e., a within-subject or time-varying covariate). In that case,  $\beta_k$  has interpretation as the change in an individual's log odds of response for a unit increase in  $X_{ijk}$ , while holding that individual's covariates fixed. Because the components of the fixed effects,  $\beta$ , have interpretations that depend on holding  $b_i$ , the  $i^{th}$  individual's random effect, fixed, they are often referred to as *subject-specific* regression coefficients. As a result generalized linear mixed models are most useful when the main scientific objective is to make inferences about individuals rather than the population averages; the population averages are the targets of inference in marginal models.

When there are between-subject (or time-invariant) covariates in the model, the interpretation of the corresponding components of  $\beta$  is somewhat less transparent and potentially misleading. If  $X_{ijk}$  is a between-subject covariate (e.g., gender, treatment, or exposure group), it is misleading to give it a subject-specific interpretation in terms of the change in the log odds of response for a unit increase in  $X_{ijk}$  since there are simply no data that provide any information about such an effect. In a sense, this interpretation of  $\beta_k$  would be a complete extrapolation beyond the observed data. Problems of interpretation with a between-subject covariate arise because a change in the value of the covariate

requires also a change in the index  $i$  of  $X_{ijk}$  to, say,  $X_{ij'k}$  (for  $i \neq i'$ ). However,  $\beta_k$  then becomes confounded with  $b_i - b_{i'}$ , the difference between the unobserved random effects for the two individuals indexed by  $i$  and  $i'$ , respectively. To circumvent this problem, we must assume that  $b_i = b_{i'}$ . That is,  $\beta_k$  must be given an interpretation in terms of a contrast of the log odds of response for two different individuals who happen to have the same value for the unobserved random effect (i.e.,  $b_i = b_{i'}$ ), but who differ by one unit in the covariate  $X_{ijk}$  (e.g., one individual is exposed, the other is unexposed, or one individual is randomized to treatment, the other to placebo). Because the random effects are latent, unobserved variables, this effect of the covariate is not directly observable from the data at hand. As a result it is somewhat unclear where the information about  $\beta_k$  is obtained when  $X_{ijk}$  is a between-subject covariate. In fact, the estimate of  $\beta_k$  is based on comparisons between subjects, not comparisons within subjects, and depends on assumptions about the distribution of the random effects. Because the estimate of  $\beta_k$  is a model-based extrapolation, it may be more sensitive to assumptions concerning the random effects distribution that are difficult to check from the data at hand.

The distinction between the regression coefficients in generalized linear mixed models and marginal models is best understood in terms of the targets of inference. In generalized linear mixed models the target of inference is the individual, since the regression coefficients have interpretation in terms of contrasts of the transformed conditional means,

$$E(Y_{ij}|X_{ij}, b_i).$$

By conditioning on the unobserved random effects,  $b_i$ , the target of inference has shifted from the population to the individual. In contrast, in marginal models the target of inference is the population, since the regression coefficients in marginal models have interpretation in terms of contrasts of the transformed population means,

$$E(Y_{ij}|X_{ij}),$$

and describe how the average response varies across different subsets of the study population defined by the covariates (e.g., gender, exposure groups, treatment groups).

Note that the population means,

$$E(Y_{ij}|X_{ij}),$$

in marginal models are averaged over the natural individual-to-individual heterogeneity in the study population (as well as over the measurement or sampling variability in the response).

For the special case where an identity link function has been adopted (i.e., for the special case of linear mixed effects models), the regression coefficients in the model for the conditional means,

$$E(Y_{ij}|X_{ij}, b_i) = X'_{ij}\beta + Z'_{ij}b_i,$$

also happen to have interpretation in terms of the population means, since

$$\begin{aligned} E(Y_{ij}|X_{ij}) &= E\{E(Y_{ij}|X_{ij}, b_i)\} \\ &= E(X'_{ij}\beta + Z'_{ij}b_i) \\ &= X'_{ij}\beta + Z'_{ij}E(b_i) \\ &= X'_{ij}\beta \end{aligned}$$

when averaged over all individuals in the study population. That is, averaged over the distribution of the random effects, the population means also follow a linear model with regression coefficients  $\beta$ . However, in general, for the non-linear link functions usually adopted for discrete data, this relationship no longer holds. That is, if

$$g\{E(Y_{ij}|X_{ij}, b_i)\} = X'_{ij}\beta + Z'_{ij}b_i,$$

where  $g(\cdot)$  is a non-linear link function (e.g.,  $\text{logit}(\cdot)$  or  $\log(\cdot)$ ), then

$$g\{E(Y_{ij}|X_{ij})\} \neq X'_{ij}\beta$$

for all  $\beta$ , when averaged over the distribution of the random effects. Thus, for non-linear (or non-identity) link functions, the regression coefficients in generalized linear mixed effects and marginal models have quite distinct interpretations, and these two classes of regression models have different

targets of inference. In short, these two classes of models address different scientific questions. Marginal models address scientific questions that are concerned with changes in the (transformed) mean response over time in the study population, and the impact of covariates on these changes. In contrast, generalized linear mixed effects models address scientific questions that are concerned with changes in the mean response for any individual, and the impact of covariates on these changes.

The following simple illustration helps highlight the main distinction between the regression coefficients in marginal and generalized linear mixed models. Consider the hypothetical data presented in [Table 14.1](#). It displays the (usually unobserved) true propensity for disease,  $\Pr(Y_{ij} = 1|b_i)$ , for three individuals measured at baseline and following treatment with a new drug intended to reduce the risk of disease. The three individuals are discernibly different in terms of their underlying propensity for disease at baseline. This heterogeneity can be expressed in terms of random effects,  $b_i$ . In a sense, individuals A, B, and C have “high,” “medium,” and “low” underlying risk for disease. Also let us assume that the entire population is composed of an equal number of individuals that fall into these three distinct risk groups. Based on this assumption, the final row of [Table 14.1](#) contains the population averages (obtained as equally weighted means).

**Table 14.1** Hypothetical data on the true propensity for disease, at baseline and post-baseline, for three individuals with heterogeneous propensities for disease.

Individual	Baseline	Post-Baseline	Difference	Log(Odds Ratio)
A	0.80	0.67	-0.13	-0.68
B	0.50	0.33	-0.17	-0.71
C	0.20	0.11	-0.09	-0.70
Population Average	0.50	0.37	-0.13	

If we considered a linear model for the probability of disease, the risk difference, or difference between the probabilities of disease at baseline and post-baseline, provides a measure of the effectiveness of the new drug. These differences (post-baseline – baseline) are displayed in the fourth column of [Table 14.1](#) and vary from -0.09 to -0.17. These can be thought of as subject-specific effects of the drug. We can then consider two possible ways to produce a single-number summary of the effectiveness of the drug. The first summary can be obtained by taking the average of the subject-specific effects (as a single-number summary of the subject-specific effects),

$$\frac{-0.13 - 0.17 - 0.09}{3} = -0.13.$$

Alternatively, the average propensity for disease at baseline (0.5) can be compared to the average propensity for disease post-baseline (0.37). The latter can be thought of as a contrast of population averages and this comparison also yields

$$(0.37 - 0.50) = -0.13.$$

That is, the difference (post-baseline – baseline) between the population averages is identical to the population average of the individual-specific differences. As such, the “difference of the averages” is equal to the “average of the differences.” This simple numerical illustration confirms the remark that was made earlier about how the fixed effects regression coefficients in the linear mixed effects model (with identity link function) also happen to have interpretation in terms of population averages.

Let us consider a non-linear function of the propensity for disease (this corresponds to adopting a non-linear link function for the probability of disease). The log odds ratio provides a natural measure of the effectiveness of the drug in reducing the risk of disease from baseline. The log odds ratios (comparing the odds of disease post-baseline to the odds of disease at baseline) for individuals A, B, and C are displayed in the fifth column of [Table 14.1](#). For example, the log odds ratio for

individual A is

$$\log \left\{ \frac{0.67/(1 - 0.67)}{0.8/(1 - 0.8)} \right\} = -0.68.$$

The log odds ratios are all very similar in magnitude, ranging from  $-0.68$  to  $-0.71$ . Once again, these can be thought of as subject-specific effects of the drug. We can then consider two possible ways to produce a single-number summary of the effectiveness of the drug. The first can be obtained by taking the average of the subject-specific effects (as a single-number summary of the subject-specific effects),

$$\frac{-0.68 - 0.71 - 0.70}{3} = -0.697.$$

This indicates that the effect of the drug on any individual is to approximately halve the odds of disease (since  $e^{-0.697} \approx 0.5$ ). Alternatively, the effectiveness of the drug can be assessed by comparing the log odds of disease in the population at baseline,  $\log(0.5/0.5) = 0$ , with the log odds of disease in the population post-baseline,  $\log(0.37/0.63) = -0.532$ . The latter can be thought of as a contrast of population log odds, and this comparison yields a measure of effect,  $-0.532$ , which is approximately 25% smaller than the summary of the subject-specific effect,  $-0.697$ . That is, the comparison of population log odds results in a discernibly different measure of the effectiveness of the drug than was found in the comparison of subject-specific effects. With a non-linear function of the propensity for disease, a “non-linear contrast of the averages” is not equal to the “average of the non-linear contrasts.” This highlights the main differences between these two approaches when a non-linear link function is adopted.

For this simple numerical illustration the reader may be curious about which of the two summary statistics,  $-0.697$  or  $-0.532$ , provides the more *realistic* estimate of the effectiveness of the drug. The answer is that they both do, although they address somewhat different scientific questions. The estimate of  $-0.697$  provides a measure of the expected change in the odds of disease for any individual treated with the drug. That is, there is an approximately 50% reduction in the odds of disease (since  $1 - e^{-0.697} \approx 0.5$ ) for any individual treated with the drug. This is the estimate that will be of most interest to an individual and his/her physician in the physician-patient context. On the other hand, the estimate of  $-0.532$  provides a measure of the expected change in prevalence of disease in the study population if everyone were to be treated with the drug. That is, there would be an expected reduction in the odds of disease in the population of approximately 40% (since  $1 - e^{-0.532} \approx 0.4$ ). This is the estimate that will be of most interest to public health researchers who are considering the potential benefits of the drug for the study population as a whole.

To provide further intuition for why the regression coefficients in generalized linear mixed models and marginal models differ, consider the following example of a logistic regression model, with normally distributed random intercepts:

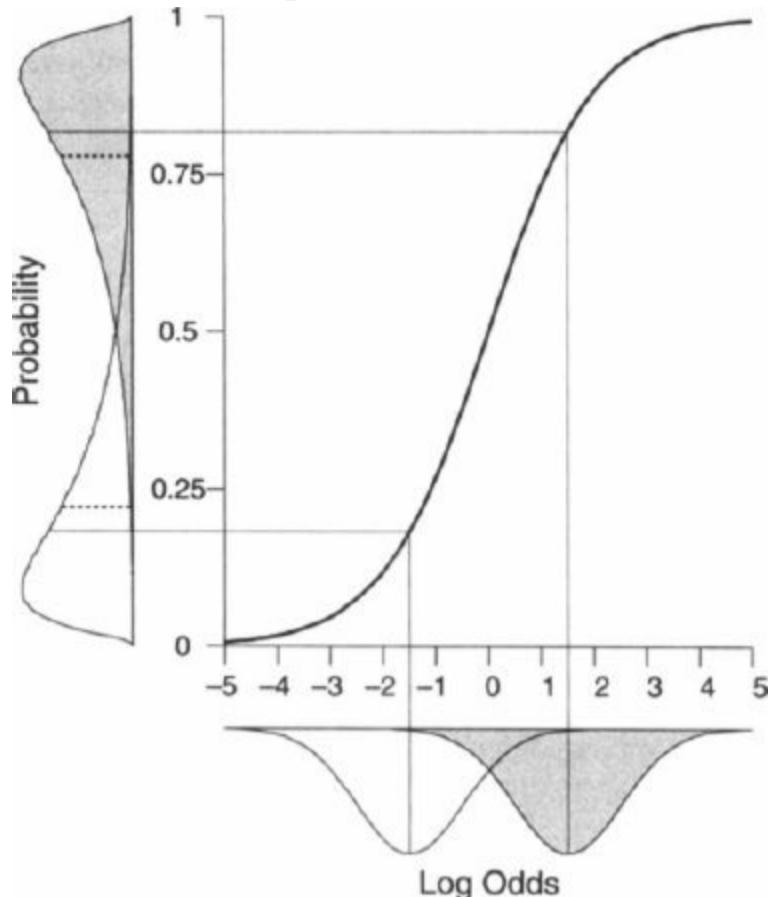
$$\text{logit}\{E(Y_{ij}|b_i)\} = \beta_1^* + \beta_2^* t_{ij} + b_i,$$

where  $t_{ij} = 0$  at baseline and  $t_{ij} = 1$  post-baseline. Similar to the illustration in [Table 14.1](#), we assume that individuals are measured at baseline and following treatment with a new drug intended to reduce the risk of disease. Individuals in the population differ in terms of their underlying propensity for disease at baseline; this heterogeneity is expressed in terms of the random effect,  $b_i$ . For a “typical” individual from the population (where a “typical” individual is one with unobserved random effect  $b_i = 0$ , the mean and median of the distribution of  $b_i$ ), the log odds of disease at baseline is  $\beta_1^*$ ; the log odds of disease following treatment with the new drug is  $\beta_1^* + \beta_2^*$ .

The log odds of disease at baseline and post-baseline are displayed in [Figure 14.1](#), for the case where  $\beta_1^* = 1.5$ ,  $\beta_2^* = -3.0$ , and  $\text{Var}(b_i) = 1.0$ . At baseline the log odds has a normal distribution with mean and median of 1.5. (See the shaded density for the log odds in [Figure 14.1](#).) From [Figure 14.1](#) it is clear that there is heterogeneity in risk of disease, with approximately 95% of individuals having a baseline log odds of disease that varies from  $-0.46$  to  $3.46$  (or  $1.5 \pm 1.96 \sqrt{1.0}$ ). When the risk of disease is translated from the log odds scale to the probability scale, the baseline probability of disease for a typical individual from the population is approximately 0.82. Furthermore

approximately 95% of individuals have a baseline probability of disease that varies from 0.39 to 0.97.

**Fig. 14.1** Subject-specific probability of disease as a function of subject-specific log odds of disease at baseline (shaded densities) and post-baseline (unshaded densities). Solid lines represent medians of the distributions; dashed lines represent means of the distributions.



From [Figure 14.1](#) it is transparent that the symmetric, normal distribution for the baseline log odds does not translate into a corresponding symmetric, normal distribution for the probability of disease. Instead, the subject-specific probabilities of disease have a negatively skewed distribution with a median, but not mean, of 0.82. (See solid line in [Figure 14.1](#).) Because of the skewness the mean of the distribution of subject-specific baseline probabilities is pulled toward the tail and is equal to 0.7785. (See dashed line in [Figure 14.1](#).) Thus the probability of disease for a “typical” individual from the population (0.82) is not the same as the prevalence of disease in the same population (0.78), due to the non-linearity of the relationship between subject-specific probabilities and log odds.

Similarly the log odds of disease post-baseline has a normal distribution with mean and median of  $-1.5$  (see the unshaded density for the log odds in [Figure 14.1](#)); approximately 95% of individuals have a post-baseline log odds of disease that varies from  $-3.46$  to  $0.46$  ( $\text{or } -1.5 \pm 1.96\sqrt{1.0}$ .) This shift in the log odds corresponds to a 20-fold decrease (since  $1 - e^{-3.0} \approx 0.95$ ) in the subject-specific odds of disease. When the risk of disease is translated from the log odds scale to the probability scale, the post-baseline probability of disease for a typical individual from the population is approximately 0.18. Furthermore approximately 95% of individuals have a post-baseline probability of disease that varies from 0.03 to 0.61. From [Figure 14.1](#) it is apparent that the subject-specific post-baseline probabilities of disease have a positively skewed distribution with median, but not mean, of 0.18. (See solid line in [Figure 14.1](#).) Because of the skewness, the mean is pulled toward the tail and is equal to 0.2215. (See dashed line in [Figure 14.1](#).)

[Figure 14.1](#) highlights how the effect of treatment on the log odds of disease for a typical individual from the population,  $\beta^*_2 = -3.0$ , is not the same as the contrast of population log odds. The latter is what is estimated in a marginal model, say

$$\text{logit}\{E(Y_{ij})\} = \beta_1 + \beta_2 t_{ij},$$

and can be obtained by comparing the log odds of disease in the population at baseline,  $\log(0.7785/0.2215) = 1.257$ , with the log odds of disease in the population post-baseline,  $\log(0.2215/0.7785) = -1.257$ . This yields a population-averaged measure of effect,  $\beta_2 = -2.514$ , which is approximately 15% smaller than  $\beta^*_2$ , the subject-specific effect of treatment.

## 14.4 OVERDISPERSION

In Chapter 11 we mentioned that overdispersion is almost the rule, not the exception, with count data. This can be potentially problematic for a mixed effects model that assumes Poisson variation for counts, conditional on the random effects; similar considerations apply to mixed effects models for binomial counts of the number of successes. Although the inclusion of random coefficients (e.g., random intercepts and slopes) induces overdispersion marginally (when averaged over the distribution of the random effects), the model nonetheless assumes Poisson variation conditional on these subject-specific effects and the covariates. One approach for relaxing the Poisson variability assumption is to extend the model to incorporate an extra source of variability in the subject-specific expected counts (or rates). For example, assuming a log-link function, the model can be extended as follows:

$$\log E(Y_{ij}|b_i, e_{ij}) = \log(T_{ij}) + X'_{ij}\beta + Z'_{ij}b_i + e_{ij},$$

where  $e_{ij}$  is an additional random effect that varies over both individuals and measurement occasions. Specifically, if a gamma distribution is assumed for the exponentiated errors,  $\exp(e_{ij})$ , with mean of 1 and variance  $\theta$ , then it can be shown that the conditional mean of  $Y_{ij}$  (conditional only on  $b_i$  and the covariates) is given by

$$\log E(Y_{ij}|b_i) = \log(T_{ij}) + X'_{ij}\beta + Z'_{ij}b_i.$$

That is, the model for the conditional mean is unchanged. However, the inclusion of the random errors implies that the corresponding conditional variance of  $Y_{ij}$  is

$$\text{Var}(Y_{ij}|b_i) = E(Y_{ij}|b_i) + \theta\{E(Y_{ij}|b_i)\}^2,$$

which is larger than the conditional mean,  $E(Y_{ij}|b_i)$ , when  $\theta > 0$  thereby allowing for overdispersion.

As was mentioned in Chapter 11, one advantage of the assumption of gamma errors is that the distribution of the counts (averaged over the distribution of these errors) has a negative binomial distribution. This makes maximum likelihood (ML) estimation of the model parameters more straightforward. That is, the model with gamma errors can be fit directly to the counts by assuming that they have a conditional negative binomial rather than a Poisson distribution, given the random effects,  $b_i$ , and the covariates.

Finally, we note that a very similar model can be specified by replacing the gamma distribution for the errors with a normal distribution. This leads to an equivalent model for the conditional variance of  $Y_i$  (conditional only on  $b_i$  and the covariates) that allows for overdispersion; also, with the exception of the intercept, the model for the conditional mean is unchanged. However, relative to a model with gamma errors, the model with normal errors is more computationally challenging to fit because the conditional distribution of  $Y_i$  (conditional only on  $b_i$  and the covariates, but averaged over the distribution of  $e_{ij}$ ) does not have a simple closed-form expression. Consequently ML estimation requires numerical integration over the distributions of random effects at two distinct levels, the level of the individual ( $b_i$ ) and the level of measurements within individuals ( $e_{ij}$ ). This is an example of a *multilevel* generalized linear mixed model, a class of models that is discussed in Chapter 22.

## 14.5 ESTIMATION AND INFERENCE

Unlike marginal models, where specification of the marginal means, variances, and pairwise associations does not fully specify the joint distribution of the vector of responses, with generalized linear mixed effects models the joint distributions of both the vector of responses and the vector of random effects are fully specified. As a result we can base estimation and inference on the likelihood function. This has important implications when there are missing data in the response variables. Specifically, conventional likelihood-based analyses of the incomplete data yield valid inferences when data are missing at random (MAR), provided that the likelihood has been correctly specified. Therefore, unlike GEE estimation of marginal models, which requires the stronger assumption that data are missing completely at random (MCAR), likelihood-based estimation of generalized linear models provides valid analyses when data are missing at random (MAR) but not MCAR. Statistical issues concerning the potential impact of missing data on analysis are discussed in detail in Chapters 17 and 18.

In this section we briefly describe maximum likelihood estimation of the fixed effects,  $\beta$ , and the random effects covariance parameters,  $G$ . We also discuss prediction of the random effects. Although ML estimation is far less straightforward for generalized linear mixed effects models than it is for the linear mixed effects models considered in Chapter 8, a variety of numerical methods for maximizing the likelihood have recently been implemented in software packages (e.g., PROC GLIMMIX in SAS, the `glmer` function in the `lme4` package in R, and the `xtmelogit` and `xtmepoisson` commands in Stata). In this section we discuss the use of quadrature methods; quadrature methods are simply numerical methods that can be made highly accurate, albeit with substantial computational overhead. In Chapter 15 we discuss two alternative methods of estimation and inference for generalized linear mixed effects models that are far less computationally demanding.

Given the three-part specification of a generalized linear mixed effects model, the joint probability for  $Y_i$  and  $b_i$  can be expressed as

$$f(Y_i, b_i) = f(Y_i|b_i)f(b_i),$$

where

$$f(Y_i|b_i) = f(Y_{i1}|b_i)f(Y_{i2}|b_i)\cdots f(Y_{in_i}|b_i)$$

under the “conditional independence” assumption. Furthermore  $f(Y_{ij}|b_i)$  is assumed to have an exponential family distribution, whereas  $f(b_i)$  is assumed to have a multivariate normal distribution, with zero mean and covariance matrix  $G$ . Since the random effects  $b_i$  are unobserved, inference about  $\beta$  and  $G$  is based on the so-called marginal or integrated likelihood function,

$$L(\beta, \phi, G) = \prod_{i=1}^N \int f(Y_i|b_i)f(b_i)db_i,$$

obtained by integrating out or averaging over the distribution of the unobserved random effects,  $b_i$ . An integral appears in the marginal likelihood, and this integral denotes the averaging over the distribution of  $b_i$ . Since the marginal likelihood has averaged over the  $b_i$ , the resulting marginal likelihood function depends only on  $\beta$ ,  $\phi$ , and  $G$ . That is, the marginal likelihood depends on the covariance of  $b_i$  but not on the unobserved  $b_i$ .

The ML estimates of  $\beta$ ,  $\phi$ , and  $G$  are simply those values of  $\beta$ ,  $\phi$ , and  $G$  that maximize this likelihood function. However, unlike the case of the linear mixed effects model, there are no simple, closed-form solutions. Instead, numerical integration techniques are required for maximizing the likelihood function. Numerical integration techniques, known as Gaussian quadrature, simply approximate the integral appearing in the marginal likelihood function as a weighted sum,

$$L(\beta, \phi, G) \approx \prod_{i=1}^N \sum_{k=1}^K f(Y_i|b_i = v_k)w_k,$$

where the known quadrature points (the weights,  $w_k$ , and the evaluation points,  $v_k$ ) are chosen to provide an accurate numerical approximation. The number of quadrature points determines the

degree of accuracy of the approximation involved in replacing the integral by a weighted sum. The number of quadrature points,  $K$ , can be increased or decreased as desired. The more quadrature points used, the more accurate the approximation will be. However, the computational burden also increases with the number of quadrature points, and grows exponentially with the number of random effects. As a result there is a trade-off that must be carefully balanced between the computational burden of quadrature methods and the desired accuracy of the results. In general, computational time is negligible when compared to the time expended in collecting longitudinal data. As a result we recommend increasing the number of quadrature points until there is evidence that all parameter estimates and standard errors are stable.

Given ML estimates of  $\beta$ ,  $\phi$ , and  $G$ , the random effects  $b_i$  for any particular subject can be predicted as follows:

$$\hat{b}_i = E(b_i|Y_i; \hat{\beta}, \hat{\phi}, \hat{G}).$$

That is, the predicted random effects for the  $i^{th}$  subject are simply “estimated” as the conditional mean of  $b_i$  given  $Y_i$  (and  $\hat{\beta}, \hat{\phi}, \hat{G}$ ); this coincides with the empirical Bayes or BLUP used for prediction of  $b_i$  in Chapter 8. Note that  $E(b_i|Y_i; \hat{\beta}, \hat{\phi}, \hat{G})$ , being a conditional mean, also requires integrating (or averaging) over the distribution of the unobserved random effects,  $b_i$ . As a result simple analytic solutions for  $b_i$  are rarely available, and numerical integration techniques must also be used here.

Finally, in our discussion of generalized linear mixed models we have assumed the distribution of the random effects is multivariate normal. Distributional assumptions about the random effects are difficult to assess from the data at hand. In particular, when the response variable is discrete, the data often contain little information to distinguish between competing distributions for the random effects. As mentioned at the end of Section 10.5, predictions of the random effects (i.e., the empirical BLUPs) are known to be heavily influenced by the normal distribution assumption for the random effects. Because the distribution of the empirical BLUPs inherits much of its shape from the assumed distribution for the random effects, histograms and normal quantile plots of the empirical BLUPs cannot be relied on for assessing the adequacy of the normal distribution assumption for the random effects. However, in general, the estimates of the fixed effects are much less sensitive to misspecification of the random effects distribution. That is, assuming that the random effects have a normal distribution when the true distribution of the random effects is non-normal (e.g., a skewed distribution) does not produce discernibly biased estimates of the fixed effects. The fixed effects estimates are, however, sensitive to a different kind of misspecification of the random effects distribution. When the assumption that the random effects are independent of the covariates,  $X_i$ , does not hold, the estimates of the fixed effects can be severely biased. This type of misspecification might arise, for example, in a study where one exposure group is more heterogeneous than another (i.e., the variance of the random effects depends on the exposure group).

## 14.6 A NOTE ON CONDITIONAL MAXIMUM LIKELIHOOD

In the previous section we described maximum likelihood estimation under the assumption that the distribution of the random effects is normal. There is an alternative approach to estimation of the fixed effects,  $\beta$ , that considers the  $b_i$  to be an additional set of fixed parameters. For example, instead of introducing randomly varying intercepts for each individual in the logistic regression model,

$$\log \left\{ \frac{\Pr(Y_{ij} = 1|b_i)}{\Pr(Y_{ij} = 0|b_i)} \right\} = X'_{ij}\beta + b_i,$$

we can incorporate fixed intercepts via the inclusion of indicator variables for each individual (in addition to the usual covariates of interest),

$$\log \left\{ \frac{\Pr(Y_{ij} = 1|\alpha_i)}{\Pr(Y_{ij} = 0|\alpha_i)} \right\} = X'_{ij}\beta + \alpha_i,$$

where the  $\alpha_i$  denote *fixed effects* representing stable (i.e., time-invariant) characteristics of individuals that are not otherwise accounted for by the covariates in the model. Although such a model is no longer a generalized linear mixed effects model (GLMM), the regression parameters,  $\beta$ , do have similar interpretations in terms of *subject-specific* effects of the covariates. This model is completely analogous to the linear *fixed effects* model described in Chapter 9. However, unlike in the linear fixed effects model, maximum likelihood estimation of the logistic regression parameters breaks down when the number of repeated measurements on each individual is relatively small. In general, optimal properties of maximum likelihood estimation require that the sample size,  $N$ , is large relative to the number of model parameters to be estimated. Herein lies the problem with maximum likelihood estimation in this setting: the number of fixed intercepts to be estimated grows as  $N$  does, while the amount of information about each parameter remains fixed. This problem is well-recognized in the statistical literature where it is referred to as the “incidental parameters problem.” Fortunately, there is a variant of maximum likelihood, known as *conditional* maximum likelihood, that can be used for inference in this setting.

The main idea behind the conditional likelihood approach is to eliminate the fixed subject-specific intercepts,  $\alpha_i$ , by constructing a likelihood that is conditional on the “sufficient statistics” for these parameters. In doing so, the “incidental parameters problem” is circumvented because the conditional likelihood no longer depends on the subject-specific intercepts. In addition the conditional likelihood has a relatively simple closed-form expression. In all other respects, though, the conditional likelihood can be used for estimation of  $\beta$  in the usual way; that is, the conditional ML estimates of  $\beta$  are those values maximizing the conditional likelihood. For example, in the logistic regression model with subject-specific intercepts, the sufficient statistic for  $\alpha_i$  is the total number of successes for each subject, say  $S_i = \sum_{j=1}^{n_i} Y_{ij}$  (for  $i = 1, \dots, N$ ). As an aside, note that conditioning on the total number of successes provides justification for the use of data on discordant pairs only in the familiar McNemar’s test for paired ( $n_i = 2$ ) binary data; in the general matched pair design, the responses for the two members of the pair are similar to repeated measures. Conditional maximum likelihood can be used for estimation of  $\beta$  in any generalized linear model with fixed subject-specific effects provided that a canonical link function is adopted and the model is restricted to having subject-specific intercepts only. We note that the fixed effects estimator of within-subject effects in the linear fixed effects model discussed in Chapter 9 (Section 9.2) can be derived as the conditional ML estimator.

There are two potential advantages of the conditional likelihood approach. First, the conditional likelihood method makes no assumptions about the distribution of the subject-specific effects. In contrast, GLMMs assume that the distribution of  $b_i$  is normal, and inferences may be sensitive to any misspecification of the random effects distribution. Moreover it is difficult to test the validity of the distributional assumption for  $b_i$  in GLMMs. Second, similar to the linear fixed effects model, the

within-subject estimator of  $\beta$  from the conditional likelihood method is not subject to confounding by omission of between-subject covariates. Therefore the conditional ML estimators are less sensitive to model misspecification than the GLMM estimators of  $\beta$ . Recall that the latter require that the random effects,  $b_i$ , be independent of the covariates,  $X_i$ . However, these two advantages are offset by the following limitations. First, the conditional likelihood approach cannot be generalized to more complex models with both subject-specific intercepts and subject-specific slopes. Second, the method can only be applied in models that assume the canonical link function (e.g., logistic regression models for Bernoulli or binomial outcomes, loglinear regression models for Poisson count data); for other (non-canonical) link functions, the “sufficient statistics” for the subject-specific effects do not exist. Third, conditional maximum likelihood does not permit estimation of the subject-specific effects or their variability. Fourth, when data are incomplete, the standard conditional ML estimator is biased if data are missing at random (MAR) but not missing completely at random (MCAR); see Section 4.3 for the definitions of, and the distinction between, MCAR and MAR. That is, the conditional ML estimator yields valid inferences when either the data are complete or any missing data are MCAR; when data are MAR, it can produce badly biased estimates of the effects of time-varying covariates. Fifth, the conditional ML estimator is usually less efficient than the GLMM estimator of  $\beta$ ; this is the price to be paid for making fewer assumptions. Finally, the conditional likelihood method can only estimate the effects of within-subject or time-varying covariates. The effects of time-invariant covariates cannot be estimated by conditional ML; this is the same limitation that was noted for linear fixed effects models in Chapter 9. This restriction to estimating only the effects of time-varying covariates is unappealing when there is scientific interest in the effects of both time-varying and time-invariant effects. Indeed, in many longitudinal studies the primary covariates of scientific interest are time-invariant, such as fixed treatment or exposure groups, or various background characteristics of individuals (e.g., gender or socioeconomic status).

In summary, conditional inference based on generalized linear models that treat the subject-specific effects as fixed rather than random can be regarded as more robust than maximum likelihood estimation under the assumptions of a corresponding GLMM. The validity of the conditional ML estimators of  $\beta$  requires fewer assumptions. On the other hand, the potential use of conditional inference for longitudinal data analysis is far more limited, restricted to the effects of within-subject covariates. Because the validity of the conditional ML estimators requires fewer assumptions, a diagnostic check on the assumptions for the random effects in GLMMs can be based on a comparison of the estimates of the within-subject effects obtained from these two approaches. Specifically, when the assumptions for the random effects in GLMMs hold, the differences between these estimates should be small; discernible differences, beyond those due to sampling variation, suggest that the assumptions regarding the random effects are not valid. This comparison can be formalized by constructing a statistical test based on the standardized differences between the two sets of estimates for the within-subject effects; this test, developed for GLMMs by Tchetgen and Coull (2006), is completely analogous to the “Hausman test” discussed in Chapter 9. We illustrate this strategy of using the conditional ML estimates as a diagnostic assessment of the assumptions for the random effects in GLMMs in Section 16.5.

## 14.7 CASE STUDIES

In this section we illustrate the main ideas presented earlier by considering GLMMs for analyzing longitudinal data from three studies. The first illustration considers a logistic regression model, with random effects, for analyzing data on amenorrhea from a randomized clinical trial of contracepting women. The second illustration considers a Poisson regression model, with random effects, for analyzing count data on epileptic seizures from a clinical trial of the anti-epileptic drug, progabide. The third example considers a proportional odds model, with random effects, for analyzing ordinal data from a clinical trial of patients with rheumatoid arthritis.

# Clinical Trial of Contracepting Women

The first example is from a longitudinal clinical trial of contracepting women reported by Machin et al. (1988). In this trial women received an injection of either 100 mg or 150 mg of depot-medroxyprogesterone acetate (DMPA) on the day of randomization and three additional injections at 90-day intervals. There was a final follow-up visit 90 days after the fourth injection, that is, one year after the first injection. Throughout the study each woman completed a menstrual diary that recorded any vaginal bleeding pattern disturbances. The diary data were used to determine whether a woman experienced amenorrhea, the absence of menstrual bleeding for a specified number of days.

A total of 1151 women completed the menstrual diaries, and the diary data were used to generate a binary sequence for each woman, according to whether she had experienced amenorrhea in the four successive three-month intervals. A feature of this clinical trial is that there was substantial dropout. More than one-third of the women dropped out before the completion of the trial; 17% dropped out after receiving only one injection of DMPA, 13% dropped out after receiving only two injections, and 7% dropped out after receiving three injections. For women who dropped out before the end of the 90-day injection interval, a determination of whether they experienced amenorrhea was made, on a proportionate basis, using their existing menstrual diary data for that interval. Statistical issues concerning the potential impact of missing data on the analysis are discussed in Chapters 17 and 18.

In clinical trials of modern hormonal contraceptives, pregnancy is exceedingly rare (and would be regarded as a failure of the contraceptive method), and this is not the main outcome of interest in this study (Machin et al., 1988). The outcome of interest is instead a binary response indicating whether a woman experienced amenorrhea in the four successive three-month intervals. The goal of the analyses presented here is to determine subject-specific changes in the risk of amenorrhea over the course of the study (12 months), and the influence of dosage of DMPA on changes in a woman's risk of amenorrhea. Of note, the treatment covariate (high versus low dosage of DMPA) is time-invariant.

Let  $Y_{ij} = 1$  if the  $i^{th}$  woman experienced amenorrhea in the  $j^{th}$  injection interval ( $j = 1, \dots, 4$ ), and  $Y_{ij} = 0$  otherwise. The following mixed effects logistic regression model for  $Y_{ij}$  was fit to the data:

$$\begin{aligned}\text{logit}\{E(Y_{ij}|b_i)\} &= \beta_1 + \beta_2 \text{Time}_{ij} + \beta_3 \text{Time}_{ij}^2 + \beta_4 \text{Dose}_i \times \text{Time}_{ij} \\ &\quad + \beta_5 \text{Dose}_i \times \text{Time}_{ij}^2 + b_i,\end{aligned}$$

where  $\text{Time} = 1, 2, 3, 4$  for the four consecutive 90-day injection intervals, and  $\text{Dose} = 1$  if randomized to 150 mg of DMPA, and  $\text{Dose} = 0$  otherwise. Note that there is no baseline measure of amenorrhea prior to receiving the first contraceptive injection. However, due to randomization, we assume that the baseline risk (at Time = 0) is the same in both dosage groups and omit a main effect of dose from the model.

Given  $b_i$ , it is assumed that the  $Y_{ij}$  are independent and have a Bernoulli distribution, with  $\text{Var}(Y_{ij}|b_i) = E(Y_{ij}|b_i) \{1 - E(Y_{ij}|b_i)\}$ , and  $\phi = 1$ . Finally, we assume that the single random effect  $b_i$  has a univariate normal distribution, with zero mean and variance  $\sigma^2_b$ ,  $b_i \sim N(0, \sigma^2_b)$ . This mixed effects model posits natural heterogeneity in women's propensity or underlying risk of amenorrhea that persists throughout all binary responses obtained over the duration of the study.

The ML estimates of the regression parameters for this model are presented in [Table 14.2](#). These results provide evidence that the subject-specific log odds of amenorrhea increases over the 12 months of the trial, and that subject-specific changes in the risk of amenorrhea depend on the dose of DMPA. For example, for a woman assigned to the low dose of DMPA, the log odds of amenorrhea increases approximately linearly, with an increase in the log odds of 1.09 (or  $1.1332 - 0.0419$ ) at 3 months, 2.10 (or  $2 \times 1.1332 - 4 \times 0.0419$ ) at 6 months, 3.02 (or  $3 \times 1.1332 - 9 \times 0.0419$ ) at 9 months, and 3.86 (or  $4 \times 1.1332 - 16 \times 0.0419$ ) at 12 months. These increases in risk correspond to odds ratios of 3.0 (or  $e^{1.09}$ ), 8.2 (or  $e^{2.10}$ ), 20.5 (or  $e^{3.02}$ ), and 47.5 (or  $e^{3.86}$ ) at 3, 6, 9, and 12 months, respectively. On the other hand, for a woman assigned to the high dose of DMPA, the log odds of amenorrhea increases quadratically, with an increase of 1.55 at 3 months, 2.79 at 6 months, 3.73 at 9 months, and 4.37 at 12 months. That is, the early linear trend shows a decline toward the

end. These increases in risk correspond to odds ratios of 4.7 (or  $e^{1.55}$ ), 16.3 (or  $e^{2.79}$ ), 41.7 (or  $e^{3.73}$ ), and 79.0 (or  $e^{4.37}$ ) at 3, 6, 9, and 12 months, respectively. For both groups, all of these increases in the odds ratios are significant at the 0.05 level.

**Table 14.2** Parameter estimates and standard errors from a mixed effects logistic regression model, with randomly varying intercepts, for the amenorrhea data.

Variable	Estimate	SE	Z
Intercept	-3.8057	0.3050	-12.48
Time	1.1332	0.2682	4.22
Time <sup>2</sup>	-0.0419	0.0548	-0.76
Dose × Time	0.5644	0.1922	2.94
Dose × Time <sup>2</sup>	-0.1096	0.0496	-2.21
$\sigma_b^2$	5.0646	0.5840	

*Note:* ML estimation is based on 50-point adaptive Gaussian quadrature.

Because treatment (low versus high dose of DMPA) is a between-subject variable, the interpretation of the fixed effects for the dose × time interactions is more difficult. The interaction effects must be given an interpretation in terms of a contrast of the increases in log odds of amenorrhea (or the odds ratio) for two different women who happen to have the same underlying risk of experiencing amenorrhea prior to randomization but who differ in terms of dose (i.e., one is assigned to low dose and the other to high dose). From the estimates of the fixed effects in Table 14.2, the ratio of the increased odds of amenorrhea at 12 months for a woman assigned to the high dose, versus another woman with the same risk of amenorrhea prior to randomization who was assigned to the low dose, is 1.66 (or  $e^{4.37-3.86}$ ) with 95% confidence interval: 1.03 to 2.66.

The estimated variance of the random intercepts is relatively large,  $\hat{\sigma}_b^2 = 5.065$ . This implies that there is substantial variability in the propensity to experience amenorrhea, since approximately 95% of the women have a baseline risk of amenorrhea that varies from

$$\frac{\exp(-3.8057 - 1.96\sqrt{5.0646})}{1 + \exp(-3.8057 - 1.96\sqrt{5.0646})}$$

to

$$\frac{\exp(-3.8057 + 1.96\sqrt{5.0646})}{1 + \exp(-3.8057 + 1.96\sqrt{5.0646})},$$

or 0.03% to 64.68%. Alternatively, we can interpret  $\hat{\sigma}_{11}$  by appealing to the notion of a latent variable distribution (see Chapter 11, Section 11.3). That is, we can assume a linear mixed effects model for the latent variable  $L_{ij}$ ,

$$L_{ij} = \beta_1 + \beta_2 \text{Time}_{ij} + \beta_3 \text{Time}_{ij}^2 + \beta_4 \text{Dose}_i \times \text{Time}_{ij} \\ + \beta_5 \text{Dose}_i \times \text{Time}_{ij}^2 + b_i + e_{ij},$$

where  $b_i$  has a normal distribution, with zero mean and variance  $\sigma_b^2$ , and the  $e_{ij}$  have a standard logistic distribution, with mean zero and variance  $\pi^2/3$ . Without loss of generality, we can assume that the threshold for categorizing  $L_{ij}$  is zero, with

$$Y_{ij} = 1 \text{ if } L_{ij} > 0,$$

$$Y_{ij} = 0 \text{ if } L_{ij} \leq 0.$$

This model for the latent variable implies the mixed effects logistic regression model for  $Y_{ij}$ ,

$$\text{logit}\{E(Y_{ij}|b_i)\} = \beta_1 + \beta_2 \text{Time}_{ij} + \beta_3 \text{Time}_{ij}^2 + \beta_4 \text{Dose}_i \times \text{Time}_{ij} \\ + \beta_5 \text{Dose}_i \times \text{Time}_{ij}^2 + b_i.$$

Thus, using the notion of an underlying latent variable distribution, we can compare the magnitudes

of the between-subject and within-subject sources of variability of the  $L_{ij}$  in terms of the intra-subject correlation (often referred to as the *intra-class correlation*)

$$\rho = \text{Corr}(L_{ij}, L_{ik}) = \frac{\sigma_b^2}{\sigma_b^2 + \pi^2/3}.$$

The estimated intra-subject correlation for the repeated latent responses is

$$\hat{\rho} = \frac{\hat{\sigma}_b^2}{\hat{\sigma}_b^2 + \pi^2/3} = \frac{5.065}{5.065 + 3.290} = 0.61,$$

indicating that there is substantial heterogeneity in the underlying propensity to experience amenorrhea. Note that  $\rho$  is the marginal correlation (averaged over the distribution of the random effects) among the unobserved  $L_{ij}$ ; it is not the marginal correlation among the  $Y_{ij}$ .

The mixed effects model considered above includes only a single random effect,  $b_i$ . With binary data, and measurements at only four occasions, greater care must be exercised in the specification of the random effects as the limited amount of data may not support estimation of more than a single variance component. Inclusion of both randomly varying intercepts and slopes for the linear time trend in the logistic regression model for the amenorrhea data resulted in convergence problems during estimation. It should not be too surprising that problems arise when fitting this model to the data. Intuitively, attempts to fit a series of logistic regressions to data on at most four observations (the number of repeated measurements on each woman with complete data; recall that due to dropout 37% of the women had fewer than four measurements) are likely to result in numerical problems and/or produce unstable estimates.

This highlights an important feature of longitudinal binary data: there is usually not much information available about random effects, beyond a random subject effect (or random intercept), when the number of repeated measurements is relatively small. Thus convergence problems during estimation are often encountered when random effects beyond a random subject effect are included in logistic regression models for longitudinal data.

The estimates of the fixed effects and variance component reported in [Table 14.2](#) were obtained by maximizing an approximate integrated likelihood, where the integration over the distribution of the random effects was achieved using numerical quadrature. Choice of the number of quadrature points determines the degree of accuracy of the approximation. In [Table 14.3](#) we display the estimate of the variance component,  $\sigma_b^2$ , and the value of the maximized log-likelihood for increasing numbers of quadrature points. The results in [Table 14.3](#) indicate that 5 to 10 quadrature points do not provide sufficient numerical accuracy for the amenorrhea data; this provides an illustration of the dangers of using too few quadrature points. The value for the maximized log-likelihood and the estimate of  $\sigma_b^2$  become stable once the number of quadrature points exceeds 30. [Table 14.3](#) also provides the CPU time required for fitting the model. As expected, the computational burden increases with the number of quadrature points. In this example there is only a single random effect and the increase in computational burden is relatively minor; however, in general, the computations grow exponentially with the number of random effects. When compared to the time expended in collecting longitudinal data, we regard the time required to accurately fit generalized linear mixed models to be negligible. So, in general, we recommend repeating analyses, with increasing number of quadrature points, until all estimates and standard errors become stable.

**Table 14.3** Sensitivity of estimate of variance component to number of quadrature points: mixed effects logistic regression model, with random intercepts, for the amenorrhea data.

Quadrature Points	Log-Likelihood	Estimate of $\sigma_b^2$	CPU Time <sup>a</sup>
1	-1957.303	4.3366	2.31
2	-1957.246	3.9811	2.49
3	-1944.100	4.3767	2.79
4	-1933.495	5.1992	3.07
5	-1936.213	4.8369	3.21
10	-1934.514	5.0539	4.45
20	-1934.465	5.0648	6.86
30	-1934.465	5.0646	9.11
40	-1934.465	5.0646	11.40
50	-1934.465	5.0646	13.73
100	-1934.465	5.0646	25.28

<sup>a</sup>CPU time (in seconds) using PROC GLIMMIX in SAS run on a PC with Intel Core 2 Duo 6600 (2.4 GHz) processor. Convergence criterion (GCONV) set equal to 1E-12.

Finally, the estimates of the fixed effects in the mixed effects logistic regression model are larger than those obtained in a similar analysis using a marginal model. For illustrative purposes we fit the following marginal model to the amenorrhea data:

$$\text{logit}\{E(Y_{ij})\} = \beta_1 + \beta_2 \text{Time}_{ij} + \beta_3 \text{Time}_{ij}^2 + \beta_4 \text{Dose}_i \times \text{Time}_{ij} \\ + \beta_5 \text{Dose}_i \times \text{Time}_{ij}^2,$$

and assume an unstructured log odds ratio pattern for the within-subject association,

$$\log \text{OR}(Y_{ij}, Y_{ik}) = \alpha_{jk},$$

where

$$\text{OR}(Y_j, Y_k) = \frac{\Pr(Y_j = 1, Y_k = 1)}{\Pr(Y_j = 1, Y_k = 0)} \frac{\Pr(Y_j = 0, Y_k = 0)}{\Pr(Y_j = 0, Y_k = 1)}.$$

The estimated regression coefficients and pairwise log odds ratios for the within-subject association, obtained using the GEE approach, are presented in [Table 14.4](#). The estimated logistic regression coefficients are smaller (in absolute value) than the estimated fixed effects in [Table 14.2](#). Furthermore, from the estimates of the fixed effects in [Table 14.4](#), the ratio of the population odds of amenorrhea at 12 months for women on the high dose versus low dose is 1.30, with 95% confidence interval: 0.98 to 1.71. These differences in the estimated coefficients and odds ratios are due to the different interpretations of  $\beta$  in the two classes of models; that is, these two classes of models estimate parameters that address substantively different questions.<sup>1</sup> The estimates of the fixed effects of dose in the mixed effects logistic regression model describe the effect of dose on a specific woman's risk of amenorrhea. The corresponding effects in the marginal logistic regression model describe the effects of dose on the prevalence of amenorrhea in the population of women assigned to high versus low doses of DMPA. Although the regression parameters for dose have distinct interpretations, their values coincide when there is no effect of dose. That is, at the null value the same hypothesis concerning the dependence of the risk of amenorrhea on dose is being tested. For example, a multivariate Wald test of  $H_0 : \beta_4 = \beta_5 = 0$  based on the marginal model parameter estimates produces  $W^2 = 12.3$ , with 2 df ( $p < 0.005$ ). The corresponding test from the mixed effects logistic regression parameter estimates produces  $W^2 = 12.4$ , with 2 df ( $p < 0.005$ ).

**Table 14.4** Parameter estimates and standard errors, obtained using GEE approach, from marginal logistic regression model for the amenorrhea data.

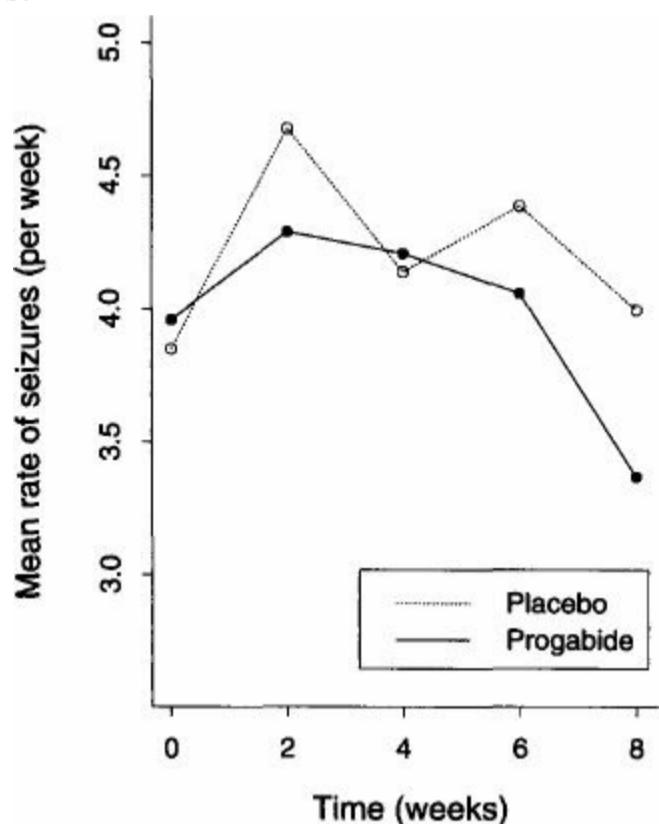
Variable	Estimate	SE	Z
Intercept	-2.2461	0.1765	-12.72
Time	0.7030	0.1581	4.45
Time <sup>2</sup>	-0.0323	0.0318	-1.02
Dose × Time	0.3380	0.1097	3.08
Dose × Time <sup>2</sup>	-0.0683	0.0284	-2.40
$\alpha_{12}$	1.8475	0.1810	10.21
$\alpha_{13}$	1.4851	0.1985	7.48
$\alpha_{14}$	1.7605	0.2482	7.09
$\alpha_{23}$	2.1610	0.1761	12.27
$\alpha_{24}$	2.0665	0.2034	10.16
$\alpha_{34}$	2.2783	0.1827	12.47

# Clinical Trial of an Anti-epileptic Drug

Next we consider data from the placebo-controlled clinical trial of 59 epileptic patients, conducted by Leppik et al. (1987). Patients with partial seizures were enrolled in a randomized clinical trial of the anti-epileptic drug, progabide. Progabide is an anti-epileptic drug whose primary mechanism of action is to enhance gamma-aminobutyric acid (GABA) content; GABA is the primary inhibitory neurotransmitter in the brain.

Participants in the study were randomized to either progabide or a placebo, as an adjuvant to the standard anti-epileptic chemotherapy. Prior to receiving treatment, baseline data on the number of epileptic seizures during the preceding 8-week interval were recorded. Counts of epileptic seizures during 2-week intervals before each of four successive post-randomization clinic visit were recorded. The average rates of seizures (per week), at baseline and in the four post-randomization visits are displayed in [Figure 14.2](#).

**Fig. 14.2** Mean rate of seizures (per week) at baseline, week 2, week 4, week 6, and week 8 in the progabide and placebo groups.



These data contain an extreme observation or outlier: one of the patients (patient 49) reported 151 seizures in the 8-week baseline interval and 302 ( $102 + 65 + 72 + 63$ ) seizures during the four successive 2-week intervals. This patient was assigned to the progabide group. Since this patient could potentially have an inordinate impact on the analysis, we present results that include and exclude data from this patient.

We consider an analysis that addresses the question of whether treatment with progabide reduces the rate of epileptic seizures (when compared to placebo). To address this question, we can compare the subject-specific changes, from baseline to follow-up, in the rate of seizures for patients in the two treatment groups. We consider the following mixed effects log-linear regression model for the subject-specific expected counts (or rates) of seizures,

$$\begin{aligned} \log E(Y_{ij}|b_i) &= \log(T_{ij}) + \beta_1 + \beta_2 \text{Time}_{ij} + \beta_3 \text{Trt}_i + \beta_4 \text{Trt}_i \times \text{Time}_{ij} \\ &\quad + b_{1i} + b_{2i} \text{Time}_{ij} \\ &= \log(T_{ij}) + (\beta_1 + b_{1i}) + (\beta_2 + b_{2i}) \text{Time}_{ij} + \beta_3 \text{Trt}_i \\ &\quad + \beta_4 \text{Trt}_i \times \text{Time}_{ij}, \end{aligned}$$

where  $Y_{ij}$  is the number of epileptic seizures for the  $i^{th}$  patient in the  $j^{th}$  period of observation ( $j = 0, \dots, 4$ ), and  $T_{13}$  is the length of period  $j$  (where  $T_{ij} = 8$  if  $j = 0$  and  $T_{ij} = 2$  if  $j = 1, 2, 3, 4$ ). The offset,  $\log(T_{ij})$ , is included because the “time at risk” is not the same in the baseline (8 weeks) and four successive follow-up periods (2-week intervals). The variable Trt is an indicator variable for

treatment group, with  $\text{Trt} = 0$  if an individual was randomized to the placebo group and  $\text{Trt} = 1$  if randomized to the progabide group. The binary variable Time denotes the baseline and follow-up periods, with  $\text{Time} = 0$  for the baseline period and  $\text{Time} = 1$  for the follow-up periods (periods 1–4). Given  $b_i$ , it is assumed that the  $Y_{ij}$  are independent and have a Poisson distribution, with  $\text{Var}(Y_{ij}|b_i) = E(Y_{ij}|b_i)$ , (i.e.,  $\phi = 1$ ). Finally, we assume that the random intercepts and slopes,  $b_i$ , have a bivariate normal distribution, with zero mean and  $2 \times 2$  covariance matrix  $G$ . This mixed effects log-linear regression model posits that there is not only natural heterogeneity among patients in terms of their baseline expected rate of seizures but also heterogeneity in the changes in the expected rates of seizures over time. Unlike the previous example, where there was not much information available in the four repeated binary response about random effects beyond a random subject effect, with these repeated count data there is sufficient information to estimate both random intercepts and slopes.

In [Table 14.5](#) we summarize the interpretation of  $\beta$  in terms of the subject-specific log expected seizure rates (per week) in the two groups at baseline and during post-baseline follow-up. Because all of the covariates in the model are dichotomous, the log-linear fixed effects regression parameters can be given interpretations in terms of subject-specific (log) rate ratios. So, for example,  $e^{\beta_2}$  is the rate ratio of seizures, comparing the follow-up periods to baseline, for a “typical” patient in the placebo group (a “typical” patient is one with unobserved random slope  $b_{2i} = 0$ , the mean and median of the distribution of  $b_{2i}$ ). Similarly  $e^{\beta_2 + \beta_4}$  is the rate ratio of seizures, comparing the follow-up periods to baseline, for a “typical” patient in the progabide group (with unobserved random slope  $b_{2i} = 0$ ). A direct comparison of the two treatments in terms of changes in the expected rates of seizures is expressible in terms of  $\beta_4$ . That is,  $\beta_4$  represents the difference between the changes in the log expected rates, comparing a patient from the progabide group to a patient from the placebo group, when the two patients are chosen so that they have the same value for the unobserved slope  $b_{2i}$ . That is,  $e^{\beta_4}$  is a ratio of rate ratios. If  $\beta_4 < 0$ , this indicates a greater reduction (or alternatively, a smaller increase) in the seizure rate from baseline for the patient assigned to the progabide group (when compared to the patient assigned to the placebo group).

**Table 14.5** Subject-specific log expected seizure rates in the two groups at baseline and during post-baseline follow-up.

Treatment Group	Period	$\log \left\{ \frac{E(Y_{ij} b_i)}{T_{ij}} \right\}$
Placebo	Baseline	$\beta_1 + b_{1i}$
	Follow-up	$(\beta_1 + b_{1i}) + (\beta_2 + b_{2i})$
Progabide	Baseline	$(\beta_1 + b_{1i}) + \beta_3$
	Follow-up	$(\beta_1 + b_{1i}) + (\beta_2 + b_{2i}) + \beta_3 + \beta_4$

For the full sample ( $N = 59$ ) the estimated fixed effects and covariance parameters from the log-linear model are displayed in [Table 14.6](#). A test of the null hypothesis,  $H_0: \beta_4 = 0$ , indicates that there is a significant time  $\times$  treatment interaction at the 0.05 level. These results suggest that there are differences between the two treatments in terms of subject-specific changes in the expected rates of seizures. In particular, there is a greater reduction in the expected seizure rate from baseline for patients treated with progabide (when compared to patients treated with the placebo). For a patient receiving the placebo, there is no expected change in the rate of seizures (or  $e^{-0.0004} \approx 1.0$ ), while for a patient treated with progabide the expected decrease in seizures is approximately 26% (or  $e^{-0.0004 - 0.3065} = e^{-0.3069} \approx 0.74$ ).

**Table 14.6** Parameter estimates and standard errors from mixed effects log-linear regression model for the seizure data.

Variable	Estimate	SE	Z
Intercept	1.0707	0.1406	7.62
Time	-0.0004	0.1097	-0.00
Trt	0.0513	0.1931	0.27
Trt × Time	-0.3065	0.1513	-2.03
$g_{11} = \text{Var}(b_{1i})$	0.5010	0.1010	
$g_{22} = \text{Var}(b_{2i})$	0.2334	0.0608	
$g_{12} = \text{Cov}(b_{1i}, b_{2i})$	0.0541	0.0559	

Note: ML estimation is based on 50-point adaptive Gaussian quadrature.

The estimated covariance parameters for the random intercepts and slopes indicate that there is substantial variability in the baseline seizure rate in the study population and also substantial variability in the patient-to-patient changes in the seizure rates in response to treatment. For example, the estimated variance of the random intercepts,  $\hat{g}_{11} = 0.501$ , implies that there is substantial patient-to-patient variability in terms of their baseline rate of seizures, since approximately 95% of the patients have a baseline seizure rate that varies from

$$\exp(1.071 - 1.96\sqrt{0.501}) \text{ to } \exp(1.071 + 1.96\sqrt{0.501}),$$

or 0.8 to 12.0 seizures per week. Similarly there is discernible heterogeneity in the patient-to-patient changes in the seizure rates. For example, approximately 95% of patients treated with progabide have changes in the rates of seizures that vary from

$$\exp(1.071 - 1.96\sqrt{0.501}) \text{ to } \exp(1.071 + 1.96\sqrt{0.501}),$$

or changes that vary from a decrease in seizures of 71% to an increase in seizures of 88%. Finally, the correlation among the random intercepts and slopes is weak, indicating that the expected change in the seizure rates is not directly related to the baseline rate of seizures.

As was noted earlier, patient 49 is an outlier with extreme counts at all occasions. While the observations on this patient are likely to inflate the variance of the random effects, especially the variance of  $b_{1i}$ , they might also have an inordinate influence on the estimates of the fixed effects parameters. To assess the impact this patient has on the results, we repeated the analysis excluding observations on this patient ( $N=58$ ). The results of this analysis are displayed in [Table 14.7](#). A test of the null hypothesis,  $H_0: \beta_4 = 0$ , indicates that there is still a significant time  $\times$  treatment interaction at the 0.05 level. These results indicate that there is a greater reduction in the expected seizure rate from baseline for patients treated with progabide (when compared to patients treated with the placebo). For a patient receiving the placebo, there is no expected change in the rate of seizures (or  $1 - e^{0.0078} \approx 0$ ), while for a patient treated with progabide the expected decrease in seizures is approximately 30% (or  $1 - e^{0.0078 - 0.3461} = 1 - e^{-0.3383} \approx 0.29$ ). Qualitatively, the results in [Table 14.7](#) are very similar to those obtained in [Table 14.6](#). As might be expected, the exclusion of patient 49 results in a noticeably smaller estimate of  $\text{Var}(b_{1i})$ .

**Table 14.7** Parameter estimates and standard errors from mixed effects log-linear regression model for the seizure data, excluding patient 49.

Variable	Estimate	SE	Z
Intercept	1.0692	0.1344	7.96
Time	0.0078	0.1070	0.07
Trt	-0.0079	0.1860	-0.04
Trt × Time	-0.3461	0.1489	-2.33
$g_{11} = \text{Var}(b_{1i})$	0.4529	0.0935	
$g_{22} = \text{Var}(b_{2i})$	0.2163	0.0587	
$g_{12} = \text{Cov}(b_{1i}, b_{2i})$	0.0151	0.0529	

*Note:* ML estimation is based on 50-point adaptive Gaussian quadrature.

Finally, in Section 14.4 we mentioned that overdispersion is almost the rule, not the exception, with count data. To allow for overdispersion, we can extend the log-linear model to incorporate an extra source of variability in the subject-specific expected counts (or rates) of seizures,

$$\begin{aligned}\log E(Y_{ij}|b_i, e_{ij}) &= \log (T_{ij}) + (\beta_1 + b_{1i}) + (\beta_2 + b_{2i}) \text{Time}_{ij} + \beta_3 \text{Trt}_i \\ &\quad + \beta_4 \text{Trt}_i \times \text{Time}_{ij} + e_{ij},\end{aligned}$$

where  $e_{ij}$  is an additional random effect that varies over both individuals and measurement occasions. Assuming that the exponentiated errors,  $\exp(e_{ij})$ , have a gamma distribution with mean of 1 and variance  $\theta$  leaves the model for the conditional mean unchanged but implies that the conditional variance of  $Y_{ij}$  is

$$\text{Var}(Y_{ij}|b_i) = E(Y_{ij}|b_i) + [E(Y_{ij}|b_i)]^2\theta.$$

This model can be fit directly to the seizure counts data by assuming that they have a negative binomial rather than a Poisson distribution, conditional on the random effects and the covariates.

To avoid any potential problems with model identification, in fitting this extended model, we set the correlation between the random intercepts and slopes to zero. This assumption is strongly supported by the results in [Table 14.7](#); for example, the likelihood ratio test of  $H_0 : g_{12} = \text{Cov}(b_{1i}, b_{2i}) = 0$  yields  $G^2 = 0.08$ , with 1 df( $p > 0.95$ ). The results of fitting the negative binomial mixed effects model to the seizure count data are summarized in [Table 14.8](#). The negative binomial model, through the inclusion of an additional random component to account for overdispersion, has led to a very discernible improvement in fit to these data. For example, the likelihood ratio test comparing the models with and without this additional random component yields  $G^2 = 77.9$ , with 1 df( $p < 0.0001$ ). A test of the null hypothesis,  $H_0: \beta_4 = 0$ , indicates that there is a significant time  $\times$  treatment interaction at the 0.05 level. In these results we see a greater reduction in the expected seizure rate from baseline for patients treated with progabide (when compared to patients treated with the placebo). For a patient receiving the placebo, there is no expected change in the rate of seizures (or  $1 - e^{0.0054} \approx 0$ ), whereas for a patient treated with progabide, the expected decrease in seizures is approximately 30% (or  $1 - e^{(0.0054 - 0.3585)} = 1 - e^{-0.3531} \approx 0.30$ ). Qualitatively, the results in [Table 14.8](#) for the fixed effects are very similar to those reported in [Table 14.7](#). What differs, though, are the estimates of the variances of the random effects. In the mixed effects Poisson model any excess variability relative to Poisson variation is partially accounted for by the variances of the random intercepts and slopes. Therefore, with the inclusion of an additional random component in the mixed effects negative binomial model to account for this excess variability (with  $\hat{\theta} = 0.1173$ ), it should not be all that surprising that the estimated variances of the random intercepts and slopes in [Table 14.8](#) are attenuated. For example, when compared to the estimated variance of the slopes in [Table 14.7](#) (where  $\theta$  is effectively assumed to be zero), the relative magnitude has been reduced by approximately 40%. By allowing for overdispersion ( $\theta > 0$ ), the estimated standard errors in [Table 14.8](#) are somewhat larger than those reported in [Table 14.7](#). The standard errors are approximately 10% larger, which is consistent with an overdispersion factor of approximately 1.2. Thus, while

there is evidence of overdispersion, the degree of overdispersion is not substantial.

**Table 14.8** Parameter estimates and standard errors from negative binomial mixed effects log-linear regression model for the seizure data, excluding patient 49.

Variable	Estimate	SE	Z
Intercept	1.0986	0.1481	7.42
Time	0.0054	0.1160	0.05
Trt	0.0033	0.2055	0.02
Trt × Time	-0.3585	0.1629	-2.20
$g_{11} = \text{Var}(b_{1i})$	0.4414	0.1040	
$g_{22} = \text{Var}(b_{2i})$	0.1278	0.0678	
$\theta = \text{Var}\{\exp(e_{ij})\}$	0.1173	0.0254	

*Note:* ML estimation is based on 50-point adaptive Gaussian quadrature.

In Section 14.4 we remarked that a very similar model can be specified by replacing the gamma distribution for the errors with a normal distribution. This leads to a completely equivalent model for the conditional variance of  $Y_i$  (conditional only on  $b_i$  and the covariates) that allows for overdispersion. For illustrative purposes we fit the extended model with normal errors,  $e_{ij}$ , to the seizure count data. We obtained remarkably similar estimates of the fixed effects and variance components to those reported in [Table 14.8](#). For example, the estimate of  $\beta_4$ , the parameter of main interest, is -0.3525 (SE = 0.1608), which is very similar to the estimate -0.3585 (SE = 0.1629) reported in [Table 14.8](#).

# Arthritis Clinical Trial

Finally, we consider longitudinal ordinal data from the longitudinal clinical trial comparing auranofin therapy (3 mg of oral gold, twice daily) and placebo for the treatment of rheumatoid arthritis (Bombardier et al., 1986). Recall, from Section 13.4, that in this six-month, randomized, double-blind trial, 303 patients with classic or definite rheumatoid arthritis were randomized to one of the two treatment groups and followed over time. The outcome variable of interest is a global impression scale (Arthritis Categorical Scale) measured at baseline (month 0), month 2, month 4, and month 6. This is a self-assessment of a patient's current arthritis, measured on a five-level ordinal scale: (1) very good, (2) good, (3) fair, (4) poor, and (5) very poor.

The goal of the analysis is to assess changes in the odds of a more favorable response over the duration of the study, and to determine whether treatment with auranofin has an influence on these changes. Letting  $Y_{ij}$  denote the ordinal response for the  $i^{th}$  subject at the  $j^{th}$  occasion, we assume that the subject-specific log odds of a more favorable response at each occasion follows the proportional odds model

$$\log \left\{ \frac{\Pr(Y_{ij} \leq k|b_i)}{\Pr(Y_{ij} > k|b_i)} \right\} = \alpha_k + \beta_1 \text{Trt}_i + \beta_2 \sqrt{\text{Month}_{ij}} \\ + \beta_3 \text{Trt}_i \times \sqrt{\text{Month}_{ij}} + b_{1i} + b_{2i} \sqrt{\text{Month}_{ij}},$$

where  $\sqrt{\text{Month}_{ij}}$  = the square-root transformation of time, in months, for the  $i^{th}$  subject at the  $j^{th}$  occasion,  $\text{Trt}_i = 1$  if the  $i^{th}$  subject is randomized to auranofin, and  $\text{Trt}_i = 0$  if randomized to placebo. This mixed effects proportional odds model allows for randomly varying intercepts and slopes for square-root transformed time. It assumes that given  $b_i$ , the  $Y_{ij}$  are independent and have a multinomial distribution. Last, we assume that the random intercepts and slopes,  $b_i$ , have a bivariate normal distribution, with zero mean and  $2 \times 2$  covariance matrix  $G$ . This model posits that there is not only heterogeneity among patients in terms of their baseline odds of a more favorable response but also heterogeneity in the changes in the odds of a more favorable response over time.

The ML estimates of the fixed effects and covariance parameters are displayed in [Table 14.9](#). A test of the null hypothesis,  $H_0: \beta_3 = 0$ , indicates that there is a significant treatment  $\times$  (square-root transformed) time interaction ( $Z = 2.80, p < 0.01$ ). These results indicate that the pattern of subject-specific changes over time in the odds of a more favorable response differs between treatment groups. In particular, relative to baseline, the odds of a more favorable response at week 6 is increased by a factor of 2.66 (or  $e^{0.4000 \times \sqrt{6}}$ ) for a patient receiving placebo but by a factor greater than 7 (or  $e^{(0.4000+0.4005) \times \sqrt{6}} = 7.11$ ) for a similar patient receiving auranofin. At the completion of the study, a patient treated with auranofin is approximately  $2 \frac{1}{2}$  times (or  $e^{0.4005 \times \sqrt{6}} = 2.67$ ) more likely to have a more favorable response when compared to a similar patient treated with placebo. As expected, due to the randomization,  $\hat{\beta}_1 \approx 0$ , indicating that patients in the two treatment groups have similar subject-specific log odds of a favorable response at baseline (or month 0).

**Table 14.9** Parameter estimates and standard errors from mixed effects proportional odds model for the auranofin clinical trial data.

Variable	Estimate	SE	Z
$\alpha_1$	-5.2945	0.3289	-16.10
$\alpha_2$	-2.0085	0.2483	-8.09
$\alpha_3$	1.0753	0.2382	4.51
$\alpha_4$	3.8683	0.3006	12.87
Trt	0.1140	0.3184	0.36
$\sqrt{\text{Month}}$	0.4000	0.1018	3.93
Trt $\times \sqrt{\text{Month}}$	0.4005	0.1432	2.80
$g_{11} = \text{Var}(b_{1i})$	4.1424	0.8549	
$g_{22} = \text{Var}(b_{2i})$	0.3243	0.1617	
$g_{12} = \text{Cov}(b_{1i}, b_{2i})$	-0.0516	0.2626	

*Note:* ML estimation is based on 30 point adaptive Gaussian quadrature.

The estimated covariance parameters for the random intercepts and slopes indicate substantial variability in the baseline odds of a more favorable response and also discernible variability in the patient-to-patient changes in the odds in response to treatment. For example, the estimated variance of the random slopes,  $\hat{g}_{22} = 0.324$ , implies patient-to-patient variability in change in the log odds over time. Approximately 95% of patients treated with placebo have slopes for time that vary over the interval

$$0.400 \pm 1.96\sqrt{0.324} = (-0.72, 1.52),$$

whereas approximately 95% of patients treated with auranofin have slopes for time that vary over the interval

$$(0.400 + 0.401) \pm 1.96\sqrt{0.324} = (-0.32, 1.92).$$

Thus, while 92% of patients treated with auranofin are expected to have positive slopes that correspond to a more favorable response, only 76% of patients treated with placebo are expected to have a more favorable response. The remaining 24% of patients treated with placebo are expected to have a less favorable response.

# 14.8 COMPUTING: FITTING GENERALIZED LINEAR MIXED MODELS USING PROC GLIMMIX IN SAS

Until recently a potential limitation of generalized linear mixed models was their computational burden. Because there is no simple closed-form solution for the marginal likelihood, numerical integration techniques are required. Maximum (marginal) likelihood estimation has only recently been implemented in standard statistical software, for example, PROC NLMIXED and PROC GLIMMIX in SAS.

To fit generalized linear mixed models, we use the GLIMMIX procedure in SAS. PROC GLIMMIX can be used to directly maximizes an approximate integrated likelihood, where the integration over the random effects is achieved using numerical quadrature. For example, to fit a logistic regression model with randomly varying intercepts to longitudinal data from two groups (coded 0 and 1), we can use the illustrative SAS commands given in [Table 14.10](#). Similarly, to fit a mixed effects log-linear regression, with randomly varying intercepts and slopes, we can use the illustrative SAS commands given in [Table 14.11](#). To fit a mixed effects log-linear regression that assumes negative binomial variability (or overdispersion relative to Poisson variability), with randomly varying intercepts and slopes, we can use the illustrative SAS commands given in [Table 14.12](#). Finally, to fit a mixed effects proportional odds regression for ordinal responses, with randomly varying intercepts and slopes, we can use the illustrative SAS commands given in [Table 14.13](#).

**Table 14.10** Illustrative commands for a mixed effects logistic regression, with randomly varying intercepts, fitted using adaptive quadrature in PROC GLIMMIX in SAS.

```
PROC GLIMMIX METHOD=QUAD(QPOINTS=50); CLASS id group;  
  MODEL y=group time group*time / DIST=BINOMIAL LINK=LOGIT S;  
  RANDOM INTERCEPT / SUBJECT=id TYPE=UN;
```

**Table 14.11** Illustrative commands for a mixed effects log-linear regression, with randomly varying intercepts and slopes, fitted using adaptive quadrature in PROC GLIMMIX in SAS.

```
PROC GLIMMIX METHOD=QUAD(QPOINTS=50);  
  CLASS id group;  
  MODEL y=group time group*time / DIST=POISSON LINK=LOG S;  
  RANDOM INTERCEPT time / SUBJECT=id TYPE=UN;
```

**Table 14.12** Illustrative commands for a mixed effects log-linear regression assuming negative binomial variance (overdispersion relative to Poisson variance), with randomly varying intercepts and slopes, fitted using adaptive quadrature in PROC GLIMMIX in SAS.

```
PROC GLIMMIX METHOD=QUAD(QPOINTS=50);  
  CLASS id group;  
  MODEL y=group time group*time / DIST=NEGBIN LINK=LOG S;  
  RANDOM INTERCEPT time / SUBJECT=id TYPE=UN;
```

**Table 14.13** Illustrative commands for a mixed effects proportional odds regression, with randomly varying intercepts and slopes, fitted using adaptive quadrature in PROC GLIMMIX in SAS.

```
PROC GLIMMIX METHOD=QUAD(QPOINTS=50);  
  CLASS id group;  
  MODEL y=group time group*time / DIST=MULT LINK=CUMLOGIT S;  
  RANDOM INTERCEPT time / SUBJECT=id TYPE=UN;
```

PROC GLIMMIX in SAS is a procedure for fitting generalized linear mixed models (among other models) using maximum likelihood (via adaptive quadrature) or using alternative methods of

estimation that are less computationally demanding. In particular, PROC GLIMMIX provides two approximate methods of estimation of the model parameters known as *penalized quasi-likelihood* (PQL) and *marginal quasi-likelihood* (MQL). In this section we focus only on maximum likelihood (ML) estimation via adaptive quadrature; at the end of Chapter 15 we discuss the PQL and MQL options in PROC GLIMMIX.

The GLIMMIX procedure is very versatile and the syntax is remarkably similar to that used in PROC MIXED. For example, PROC GLIMMIX has a RANDOM statement that is used in a similar way as in PROC MIXED for introducing random effects. No attempt is made here to give a comprehensive review of the main features of PROC GLIMMIX. Instead, we present illustrative commands for fitting generalized linear mixed effects models in general terms (see [Tables 14.10](#) through [14.13](#)) and then describe the most salient parts of the command syntax for these illustrations. Next we present a brief description of each of the command statements used in [Tables 14.10](#) through [14.13](#).

PROC GLIMMIX <options>;

This statement calls the procedure GLIMMIX in SAS. It includes an option for specifying the method of estimation, using METHOD=<options>.

The default is METHOD=RSPL, where RSPL denotes a version of *penalized quasi-likelihood* (PQL) estimation. Penalized quasi-likelihood estimation is discussed in Chapter 15. PROC GLIMMIX also includes an implementation of adaptive Gaussian quadrature using the METHOD=QUAD option. When the METHOD=QUAD option is used, the accuracy of the numerical approximation can be increased by specifying the number of quadrature points used during evaluation of integrals for the marginal likelihood. For example, METHOD=QUAD(QPOINTS=50) specifies that 50 quadrature points be used for each random effect, resulting in a total of  $50^q$  quadrature points (where  $q$  is the number of random effects, the dimension of bi). A note of caution, the likelihood approximation may not be accurate if too few quadrature points are used.

CLASS variables;

The CLASS statement is used to define all variables that are to be regarded as categorical or factors. By default, this statement will create indicator variables for each factor using a reference group coding, with the last level (where “last” here refers to the level with the largest alphanumeric value) regarded as the reference group. Different sort orders for the CLASS variables can be requested by the ORDER=<option> on the PROC GLIMMIX statement.

MODEL response = <effects> / <options>;

MODEL events/trials = <effects> / <options>;

The MODEL statement specifies the response variable and the covariate effects. The second form of the MODEL statement, with the events/trials syntax, allows the response to be in the form of a ratio of two variables (e.g., counts of the number of successes and the number of trials) and is used for binomial response data. The linear predictor can include both discrete (defined in the CLASS statement) and quantitative (excluded from the CLASS statement) covariates. By default, PROC GLIMMIX includes a column of 1's for the intercept in the model. The SOLUTION or S option is used to produce estimates of the fixed (or covariate) effects.

The option DIST=*keyword* specifies the conditional distribution of the response given the random effects in a GLMM.

The LINK=*keyword* specifies the choice of built-in link function relating the mean response to the linear predictor. If the LINK=*keyword* is omitted, the default link function is the canonical link function associated with the particular exponential family distribution specified on DIST=*keyword*.

A final option often required when modeling count data is an offset. The OFFSET=*variable* specifies a variable to be used as an offset. For example, in modeling count data, the rate is often of more direct interest and the denominator for the counts or “population at risk” (or more specifically, the log of the denominator) can be included as an offset. Note that this variable

cannot be a CLASS variable and that it should not be included as one of the covariates listed on the MODEL statement.

**RANDOM <random-effects> / SUBJECT=subject-effect <options>;**

In a generalized linear mixed model, the RANDOM statement is used to define the covariates in the design matrix,  $Z_i$ , for the random effects,  $b_i$ . Ordinarily these will be a subset of the covariates included on the MODEL statement. While the MODEL statement is used to define the design matrix for the fixed effects and the RANDOM statement is used to define the design matrix for the random effects, note that an intercept is included by default in the former but not the latter. That is, unlike the MODEL statement, PROC GLIMMIX does not include an intercept in the RANDOM statement by default. However, you can specify INTERCEPT (or INT) as a random effect on the RANDOM statement. The SUBJECT option on the RANDOM statement is used to denote a variable that distinguishes clusters of correlated responses. By including a variable in the SUBJECT option (e.g., a subject identifier), pairs of observations with distinct values of that variable are regarded as independent. Pairs of observations with the same values of that variable share common values of the random effects. The RANDOM statement is also used to specify the structure of the covariance matrix for the random effects, G. The structure of G is specified using the TYPE=option. The random effects can be assumed to be correlated (TYPE=UN) or uncorrelated (TYPE=VC); ordinarily, covariance pattern models are not used to account for the covariance among the random effects. Finally, the SOLUTION option on the RANDOM statement produces the empirical Bayes predictions of the random effects.

We conclude by noting that PROC GLIMMIX can fit a broader class of generalized linear mixed models than have been considered in this chapter. To do so, PROC GLIMMIX makes a distinction between two sources of variation and covariation in the model for the data: (1) variation due to random effects,  $b_i$ , and (2) “residual” (co)variation. To distinguish these two sources of (co)variation, PROC GLIMMIX refers to (1) as “G-side” effects and (2) as “R-side” effects. These two non-standard terms are derived from syntax for the covariance matrices for the random effects (denoted G) and the “residual errors” (denoted R) in PROC MIXED for linear mixed effects models. PROC GLIMMIX is quite versatile in allowing models to be fit with various combinations of “G-side” and “R-side” effects. For example, in a standard generalized linear mixed model, the introduction of random effects is handled by including “G-side” effects and the “conditional independence” assumption is handled by implicitly assuming a simple structure for the “R-side” effects (i.e., uncorrelated residual errors, the default option for the “R-side” effects). In principle, it is possible to fit models that relax the “conditional independence” assumption by allowing for a more general structure for the “R-side” effects (e.g., autoregressive residual errors); however, we caution that as with linear mixed effects models, there can be subtle issues of model identification when a more general structure for the “R-side” effects is assumed because it may not be possible to estimate both the “G-side” and the more general “R-side” effects from the data at hand. A more detailed description of “G-side” and “R-side” effects can be found at the end of Chapter 15.

Finally, a word of caution concerning the use of PROC GLIMMIX. Our limited experience with this procedure indicates that it can be sensitive to poor choices of starting values for the covariance parameters for the random effects and/or the numerical accuracy of the quadrature used. Convergence of the algorithm implemented in PROC GLIMMIX should never be taken for granted; neither should convergence to a global maximum be assumed. Instead, we recommend that users of this procedure provide different initial values for the covariance parameters and/or consider a grid search of values to ensure that a global maximum has been obtained. Accurate initial values for the covariance parameters can be obtained by specifying a grid of feasible values using the PARMS statement in PROC GLIMMIX. The PARMS statement specifies initial values for the covariance or scale parameters. If you specify more than one set of initial values on the PARMS statement, PROC GLIMMIX will first evaluate the marginal likelihood at each grid value and select the grid point that produces the largest value of the marginal likelihood as the initial values for the covariance parameters for the subsequent maximization of the marginal likelihood. Problems with convergence and computation are likely to arise when the model parameters are on scales that vary by more than a

few orders of magnitude. The latter problem can be circumvented by appropriately rescaling each parameter. For example, when the variances of random intercepts and slopes differ by more than a few orders of magnitude, the variance of the slopes can be rescaled by multiplying the variable for time by an appropriate constant. The numerical accuracy of the quadrature used can always be enhanced by increasing the number of quadrature points; of course, increasing the number of quadrature points does raise the computational burden.

## 14.9 FURTHER READING

A relatively non-technical discussion of random effects models for binary data can be found in Section 6.7 of Collett (1991). Chapter 12 of Agresti (2002) provides a detailed, although more mathematically challenging, description of generalized linear mixed effects models for categorical data.

# Bibliographic Notes

The theoretical foundation for generalized linear mixed effects models can be found in Skellam's (1948) introduction of the beta-binomial distribution. Since then, the statistical literature on generalized linear mixed effects models has grown rapidly. Some key references in the literature include Cox (1970), Pierce and Sands (1975), Williams (1982), Stiratelli et al. (1984), Anderson and Aitkin (1985), Gilmour et al. (1985), Wong and Mason (1985), Schall (1991), Zeger and Karim (1991), Breslow and Clayton (1993), and Hedeker and Gibbons (1994). Molenberghs et al. (2010) present a broad class of generalized linear models that accommodate both overdispersion and correlation via the introduction of two separate sets of random effects.

The marginal maximum likelihood method described previously is based on a numerically integrated likelihood function and requires the computation of the integral over the random effects. A method known as adaptive Gaussian quadrature is commonly used for computing this integral and is described in detail in Pinheiro and Bates (1995); also see Anderson and Aitkin (1985) and Hedeker and Gibbons (1994, 1996). Lesaffre and Spiessens (2001) provide a striking example of the dangers of using too few quadrature points when fitting generalized linear mixed effects models; also see Problem 14.1.10. An alternative approximation, leading to an approach known as penalized quasi-likelihood (PQL), was proposed by Stiratelli et al. (1984). A more accurate approximation, based on higher-order Laplace approximations, is described in Breslow and Lin (1995) and Lin and Breslow (1996).

Finally, Neyman and Scott (1948) defined the so-called incidental parameters problem and showed that maximum likelihood estimation (MLE) may be problematic when the number of model parameters grows with the number of observations. The use of conditional likelihoods, and properties of conditional ML estimators, is discussed in Anderson (1970). In the longitudinal setting, Conaway (1992) and Rathouz (2004) discuss conditional ML estimation with incomplete response data.

## Problems

**14.1** In a randomized, double-blind, parallel-group, multicenter study comparing two oral anti-fungal treatments (200 mg/day Itraconazole and 250 mg/day Terbinafine) for toenail infection (De Backer et al., 1998; also see Lesaffre and Spiessens, 2001), patients were evaluated for the degree of onycholysis (the degree of separation of the nail plate from the nail-bed) at baseline (week 0) and at weeks 4, 8, 12, 24, 36, and 48 thereafter. The onycholysis outcome variable is binary (none or mild versus moderate or severe). The binary outcome was evaluated on 294 patients comprising a total of 1908 measurements. The main objective of the analyses is to compare the effects of the two oral anti-fungal treatments (Itraconazole and Terbinafine) on changes in the probability of the binary onycholysis outcome over the duration of the study.

The raw data are stored in an external file: `toenail.dat`

Each row of the data set contains the following five variables:

ID Y Treatment Month Visit

*Note:* The binary onycholysis outcome variable  $Y$  is coded 0 = none or mild, 1 = moderate or severe. The categorical variable Treatment is coded 1 = Terbinafine, 0 = Itraconazole. The variable Month denotes the exact timing of measurements in months. The variable Visit denotes the visit number (visit numbers 1–7 correspond to scheduled visits at 0, 4, 8, 12, 24, 36, and 48 weeks).

**14.1.1** First, consider a *marginal* model for the log odds of moderate or severe onycholysis. Using GEE, fit a model that assumes linear trends for the log odds over time, with common intercept for the two treatment groups, but different slopes:

$$\text{logit}\{E(Y_{ij})\} = \beta_1 + \beta_2 \text{Month}_{ij} + \beta_3 \text{Treatment}_i \times \text{Month}_{ij}.$$

Assume “exchangeable” log odds ratios (or “exchangeable” correlations, if available software does not permit the within-subject association to be parameterized in terms of log odds ratios) for the association among the repeated binary responses.

**14.1.2** What is the interpretation of  $\beta_2$  in this model?

**14.1.3** What is the interpretation of  $\beta_3$  in this model?

**14.1.4** From the results of the analysis for Problem 14.1.1, what conclusions do you draw about the effect of treatment on changes in the log odds of moderate or severe onycholysis over time? Provide results that support your conclusions.

**14.1.5** Next consider a generalized linear mixed model, with randomly varying intercepts, for the patient-specific log odds of moderate or severe onycholysis. Using maximum likelihood (ML), fit a model with linear trends for the log odds over time and allow the slopes to depend on treatment group:

$$\text{logit}\{E(Y_{ij}|b_i)\} = (\beta_1 + b_i) + \beta_2 \text{Month}_{ij} + \beta_3 \text{Treatment}_i \times \text{Month}_{ij},$$

where, given  $b_i$ ,  $Y_{ij}$  is assumed to have a Bernoulli distribution. Assume that  $b_i \sim N(0, \sigma^2_b)$ .

**14.1.6** What is the estimate of  $\sigma^2_b$ ? Give an interpretation to the magnitude of the estimated variance?

**14.1.7** What is the interpretation of the estimate of  $\beta_2$ ?

**14.1.8** What is the interpretation of the estimate of  $\beta_3$ ?

**14.1.9** Compare and contrast the estimates of  $\beta_3$  from the marginal and mixed effects models. Why might they differ?

**14.1.10** Repeat the analysis from Problem 14.1.5 sequentially increasing the number of quadrature points used. Compare the estimates and standard errors of the model parameters when the number of quadrature points is 2, 5, 10, 20, 30, and 50. Do the results depend on the number of quadrature points?

**14.2** The Skin Cancer Prevention Study was a randomized, double-blind, placebo-controlled clinical trial of beta carotene to prevent non-melanoma skin cancer in high-risk subjects (Greenberg et al., 1989, 1990; also see Stukel, 1993). A total of 1805 subjects were randomized to either placebo or 50 mg of beta carotene per day for 5 years. Subjects were examined once a year and biopsied if a cancer was suspected to determine the number of new skin cancers occurring since the last exam. The outcome variable is a count of the number of new skin cancers per year. The outcome was evaluated on 1683 subjects comprising a total of 7081 measurements. The main objective of the analyses is to compare the effects of beta carotene on skin cancer rates.

The raw data are stored in an external file: `skin.dat`

Each row of the data set contains the following 9 variables:

ID Center Age Skin Gender Exposure Y Treatment Year

*Note:* The outcome variable  $Y$  is a count of the number of new skin cancers per year. The categorical variable Treatment is coded 1 = beta carotene, 0 = placebo. The variable Year denotes the year of follow-up. The categorical variable Gender is coded 1 = male, 0 = female. The categorical variable Skin denotes skin type and is coded 1 = burns, 0 = otherwise. The variable Exposure is a count of the number of previous skin cancers. The variable Age is the age (in years) of each subject at randomization.

**14.2.1** Consider a generalized linear mixed model, with randomly varying intercepts, for the subject-specific log rate of skin cancers.

Using maximum likelihood (ML), fit a model with linear trends for the log rate over time and allow the slopes to depend on treatment group:

$$\text{log}\{E(Y_{ij}|b_i)\} = (\beta_1 + b_i) + \beta_2 \text{Year}_{ij} + \beta_3 \text{Treatment}_i \times \text{Year}_{ij},$$

where, given  $b_i$ ,  $Y_{ij}$  is assumed to have a Poisson distribution. Assume that  $b_i \sim N(0, \sigma^2_b)$ .

**14.2.2** What is the estimate of  $\sigma^2_b$ ? Give an interpretation to the magnitude of the estimated variance?

**14.2.3** What is the interpretation of the estimate of  $\beta_2$ ?

**14.2.4** What is the interpretation of the estimate of  $\beta_3$ ?

**14.2.5** From the results of the analysis for Problem 14.2.1, what conclusions do you draw about the effect of beta carotene on the log rate of skin cancers? Provide results that support your conclusions.

**14.2.6** Obtain the predicted (empirical BLUP) random effect for each subject.

(a) Calculate the sample variance of the predictions. How does it compare to the estimate of  $\sigma^2_b$  obtained in Problem 14.2.2? Why might they differ?

(b) Plot the predictions against age and the count of the number of previous skin cancers. What do you conclude?

**14.2.7** Repeat the analysis from Problem 14.2.1 adjusting for skin type, age, and the count of the number of previous skin cancers. What conclusions do you draw about the effect of beta carotene on the adjusted log rate of skin cancers?

**14.3** In the U.S. National Institute of Mental Health (NIMH) Schizophrenia Collaborative Study, 437 patients were randomly assigned to receive one of four medications: placebo, chlorpromazine, fluphenazine, or thioridazine; the latter three medications are anti-psychotic drugs. The study protocol called for longitudinal measurements to be made at weeks 0, 1, 3, and 6. A very small subset of patients were additionally measured at weeks 2, 4, and 5. The outcome variable of interest is a 4-level ordinal scale measuring “severity of illness,” derived from item 79 of the Inpatient Multidimensional Psychiatric Scale (Lorr and Klett, 1966). The four categories of the ordinal scale correspond to: 1 = “normal or borderline mentally ill,” 2 = “mildly or moderately ill,” 3 = “markedly ill,” and 4 = “severely or among the most extremely ill” (Gibbons and Hedeker, 1994; also see Hedeker and Gibbons, 2006). In this study there was substantial attrition, especially in the placebo group; however, for the purpose of this exercise, you can ignore the potential impact of missing data on the analysis. The main objective of the analyses is to assess changes in the odds of a more favorable response over the duration of the study, and to determine whether treatment with anti-psychotic drugs has an influence of these changes.

The raw data are stored in an external file: `schizophrenia.dat`

Each row of the data set contains the following four variables:

ID Y Trt Week

*Note:* The ordinal outcome variable  $Y$  is coded 1 = “normal or borderline mentally ill,” 2 = “mildly or moderately ill,” 3 = “markedly ill,” and 4 = “severely or among the most extremely ill.” The categorical treatment variable  $Trt$  is coded 1 = anti-psychotic drug, 0 = placebo; that is, the three anti-psychotic medications (chlorpromazine, fluphenazine, or thioridazine) are combined into a single group to be compared to placebo. The variable  $Week$  denotes the timing of measurements in weeks.

**14.3.1** Consider a proportional odds model, with randomly varying intercepts and slopes, for the patient-specific cumulative log odds of response.

Using maximum likelihood (ML), fit a model with linear trends for the cumulative log odds over time and allow the slopes to depend on treatment group:

$$\log \left\{ \frac{\Pr(Y_{ij} \leq k|b_i)}{\Pr(Y_{ij} > k|b_i)} \right\} = \alpha_k + \beta_1 Trt_i + \beta_2 Week_{ij} \\ + \beta_3 Trt_i \times Week_{ij} + b_{1i} + b_{2i} Week_{ij},$$

where, given  $b_{1i}$  and  $b_{2i}$ ,  $Y_{ij}$  is assumed to have a multinomial distribution. Assume that  $b_{1i}$  and  $b_{2i}$  have a bivariate normal distribution, with zero mean and  $2 \times 2$  covariance matrix  $G$ .

**14.3.2** What is the interpretation of the estimate of  $\beta_2$ ?

**14.3.3** What is the interpretation of the estimate of  $\beta_3$ ?

**14.3.4** From the results of the analysis from Problem 14.3.1, what conclusions do you draw about the effect of treatment with anti-psychotic medications on changes in the log odds of a more favorable response over time? Provide results that support your conclusions.

**14.3.5** Repeat the analysis from Problem 14.3.1 replacing time (in weeks) with square-root

transformed time:

$$\log \left\{ \frac{\Pr(Y_{ij} \leq k|b_i)}{\Pr(Y_{ij} > k|b_i)} \right\} = \alpha_k + \beta_1 \text{Trt}_i + \beta_2 \sqrt{\text{Week}_{ij}} \\ + \beta_3 \text{Trt}_i \times \sqrt{\text{Week}_{ij}} + b_{1i} + b_{2i} \sqrt{\text{Week}_{ij}}.$$

**14.3.6** Based on the results from Problem 14.3.5, estimate the change in the odds of a more favorable response at week 6, relative to baseline, for a patient receiving placebo.

**14.3.7** Based on the results from Problem 14.3.5, estimate the change in the odds of a more favorable response at week 6, relative to baseline, for a patient receiving an anti-psychotic medication.

**14.3.8** Based on the results from Problem 14.3.5, approximately what percentage of patients receiving placebo are expected to show improvement over time?

**14.3.9** Based on the results from Problem 14.3.5, approximately what percentage of patients receiving anti-psychotic medications are expected to show improvement over time?

**14.3.10** From the results of the analysis for Problem 14.3.5, what conclusions do you draw about the effect of treatment with anti-psychotic medications on changes in the log odds of a more favorable response over time? Provide results that support your conclusions.

**14.3.11** In Problem 14.3.1 it is assumed that changes in the cumulative log odds are linear in time (in weeks). In Problem 14.3.5 it is assumed that changes in the cumulative log odds are linear in square-root transformed time ( $\sqrt{\text{weeks}}$ ). Which of the two models appears to fit the data better? Provide results that support your conclusions.

**14.3.12** Repeat the analysis from Problem 14.3.5 sequentially increasing the number of quadrature points used. Compare the estimates and standard errors of the model parameters when the number of quadrature points is 2, 5, 10, 20, 30, and 50. Do the results depend on the number of quadrature points?

# *Chapter 15*

## *Generalized Linear Mixed Effects Models: Approximate Methods of Estimation*

### **15.1 INTRODUCTION**

In the previous chapter we reviewed the broad class of regression models known as *generalized linear mixed effects models* (GLMMs) and described how they extend the conceptual approach of the linear mixed effects models to generalized linear models. By allowing a subset of the regression parameters in generalized linear models to vary from individual to individual, via the introduction of random effects, GLMMs account for natural heterogeneity in the study population. However, although GLMMs are conceptually straightforward to develop, model fitting and parameter estimation can be challenging because the likelihood to be maximized does not have a simple closed-form expression. Because the likelihood to be maximized is obtained by integrating or averaging over the distribution of the random effects, maximizing that likelihood is challenging.

In Chapter 14 we discussed ML estimation of the fixed effects and variance components in GLMMs through the use of adaptive Gaussian quadrature. This approach to estimation and inference can be considered “true” ML to the extent that an adequate number of quadrature points have been used to ensure good numerical accuracy. Recall that Gaussian quadrature approximates the likelihood to be maximized by replacing the integrals in the likelihood with summations. This approximation can be made very accurate by increasing the number of quadrature points. Of course, increasing the number of quadrature points raises the computational burden. For example, a tripling of the number of quadrature points for a GLMM with two random effects will lead to an almost 10-fold increase in computation.

Methods of estimation and inference that rely on numerical quadrature can become computationally burdensome as the number of random effects in the GLMM increases. For example, while adaptive Gaussian quadrature is relatively straightforward in a GLMM with only two random effects, it can be quite hopeless in the setting of a GLMM with more than 10 random effects where it breaks down due to the “curse of dimensionality.” To see why, note that 30-point quadrature in a GLMM with two random effects requires numerical evaluations at a total of 900 (or  $30^2$ ) quadrature points. In contrast, 30-point quadrature in a GLMM with 10 random effects requires numerical evaluations at over 590 trillion (or  $30^{10}$ ) quadrature points! The computations grow exponentially with the number of random effects. This has provided the impetus for the development of alternative approximations that are far less computationally demanding. These alternative methods of estimation and inference for GLMMs are the topic of this chapter. Specifically, we review two approximate methods of estimation known as *penalized quasi-likelihood* (PQL) and *marginal quasi-likelihood* (MQL). Both of these approximate methods have been implemented in various statistical software packages (e.g., PROC GLIMMIX in SAS) and stand-alone programs that have been specifically tailored for fitting GLMMs (e.g., MLwiN). However, as we will see later, the PQL method is the only *legitimate* approximate method in the sense of yielding estimates of the parameters of the GLMM. In contrast, the MQL method, although derived as an approximate method of estimation for GLMMs, does not actually yield estimates of the parameters of the GLMM. This critical distinction between the PQL and MQL methods is one that is not well recognized; in Section 15.4 we highlight the important differences between the methods and their consequences for inference.

Finally, we emphasize at the outset that in many longitudinal studies there is only a single level of nesting in the data (i.e., repeated occasions nested within individuals) and the dimension of the

vector of random effects (e.g.,  $q$ ) is relatively small (e.g.,  $q \leq 3$ ). Moreover, with longitudinal binary data, there is usually not much information available about random effects beyond a single variance component, especially when the number of repeated measurements is relatively small. Thus in the longitudinal setting the need for these alternative approximate methods is somewhat less acute than in other settings with cluster-correlated data where there may be numerous levels of nesting in the data and where the dimensionality of the random effects may be relatively large. This is important because many of the alternative approximate methods discussed in this chapter can perform poorly in certain settings relative to adaptive Gaussian quadrature. In addition the use of PQL and MQL methods is even more problematic when longitudinal data are incomplete. The inaccuracy of their approximations implies that the methods are not valid when data are missing at random (MAR), but not missing completely at random (MCAR); see Section 4.3 for the definitions of, and the distinction between, MCAR and MAR. As a result these methods can produce badly biased estimates of the effects of covariates when data are MAR, but not MCAR. Therefore we consider ML estimation based on adaptive Gaussian quadrature, with a sufficiently large number of quadrature points, to be the method of choice for fitting GLMMs to longitudinal data.

## 15.2 PENALIZED QUASI-LIKELIHOOD

Penalized quasi-likelihood (PQL) estimation of GLMMs has been proposed as a computationally simple alternative to methods based on numerical quadrature, especially when the number of random effects is relatively large. There are two closely related derivations of the PQL algorithm. The first is obtained by applying an approximation to the integrand (i.e., the function to be integrated or averaged over) in the GLMM likelihood so that a closed-form expression for the integral is achieved. The particular approximation used is known as a conventional (or lower-order) Laplace approximation. Interestingly there is a connection between this Laplace approximation and Gaussian quadrature. It can be shown that a Laplace approximation to the integrand in the likelihood for a GLMM corresponds to the use of adaptive Gaussian quadrature with only a single quadrature point. The Laplace approximation is accurate in cases where the conditional distribution of the vector of responses, given the random effects, is approximately normal but can provide a poor approximation otherwise. Later we discuss settings where the approximation is likely to be accurate versus settings where it may yield badly biased estimates. The use of a conventional Laplace approximation is one motivation for PQL.

The second, and closely related, derivation for PQL follows from an approximation to the GLMM. Recall from Chapter 14 that the GLMM can be expressed as

$$(15.1) \quad g\{E(Y_{ij}|b_i)\} = \eta_{ij}^b = X'_{ij}\beta + Z'_{ij}b_i$$

for some known link function,  $g(\cdot)$ , and a  $q \times 1$  vector of random effects  $b_i \sim N(0, G)$ . In a very slight departure from the notation used in Chapter 14, the superscript “b” in  $\eta_{ij}^b$  (and, later, in  $\mu_{ij}^b$ ) is used to emphasize conditioning on the random effects,  $b_i$ . Because the conditional distribution of each  $Y_{ij}$ , given  $b_i$ , belongs to the exponential family of distributions (e.g., Bernoulli or Poisson),  $\text{Var}(Y_{ij}|b_i) = v\{E(Y_{ij}|b_i)\} \phi$ , where  $v(\cdot)$  is a known variance function. Note that the model for the conditional mean given by (15.1) can also be expressed as:

$$(15.2) \quad Y_{ij} = g^{-1}(X'_{ij}\beta + Z'_{ij}b_i) + \epsilon_{ij} = \mu_{ij}^b + \epsilon_{ij},$$

where  $g^{-1}(\cdot)$  is the *inverse* link function,  $\mu_{ij}^b$  denotes the *conditional* mean of  $Y_{ij}$ , given  $b_i$ , and the errors  $\epsilon_{ij}$  are assumed to have a mean of zero. The PQL method proceeds by approximating the model given by (15.2), so that a linear mixed effects model holds, at least approximately, for a transformation of the response (and a related transformation of the errors); the transformation of the response is denoted by  $Y_{ij}^*$ . Specifically, the standard linear mixed effects model,

$$(15.3) \quad Y_{ij}^* \approx X'_{ij}\beta + Z'_{ij}b_i + \epsilon_{ij}^*,$$

holds approximately for the “working” response (or so-called pseudo-data),  $Y_{ij}^*$ . The “working” response  $Y_{ij}^*$ , a transformation of  $Y_{ij}$ , depends on current estimates of the fixed effects,  $\beta$ , and the random effects,  $b_i$ ; similarly the transformed error,  $\epsilon_{ij}^*$ , depends on current estimates of  $\beta$  and  $b_i$ . Specifically, when a canonical link function has been chosen,

$$Y_{ij}^* = v^{-1}(\hat{\mu}_{ij}^b)(Y_{ij} - \hat{\mu}_{ij}^b) + X'_{ij}\hat{\beta} + Z'_{ij}\hat{b}_i$$

and  $\epsilon_{ij}^* = v^{-1}(\hat{\mu}_{ij}^b)\epsilon_{ij}$ . Recall from Section 11.2 that canonical link functions are simply unique transformations of the mean that can be derived for any selected distribution (e.g., logit link function for a Bernoulli response or log link function for a Poisson response).

Although the technical details behind this approximation are not important, we note that (15.3) is referred to as a first-order Taylor series expansion of (15.2) around current estimates  $\hat{\beta}$  and  $\hat{b}_i$ . A very similar type of expression for  $Y_{ij}^*$  can be derived when a non-canonical link function is adopted. A brief derivation of the approximation, together with an illustration of the form of the “working” response  $Y_{ij}^*$  for a mixed effects logistic regression model, can be found in Section 15.7; readers who find the level of detail in this section challenging can omit Section 15.7 at first reading without loss of continuity.

Recognizing that the “working” response  $Y_{ij}^*$  follows a linear mixed effects model, with fixed

effects  $\beta$ , and with random effects  $b_i$  and within-subject errors  $\epsilon_{ij}^*$ , estimation can proceed by iterating between the following two steps:

1. Fit a linear mixed effects model to the “working” response  $Y_{ij}^*$ , to obtain updated estimates of  $\beta$  and  $G$ , and subsequently, empirical BLUP predictions of  $b_i$ . The linear mixed effects model is fitted using weights,  $w_{ij}$ , that are inversely proportional to the variance of the  $\epsilon_{ij}^*$ ; for example,  $w_{ij} = \phi^{-1} v(\hat{\mu}_{ij}^b)$  when a canonical link function is adopted.
2. Use the estimates of  $\beta$  and  $b_i$  obtained in step 1 to update the “working” response  $Y_{ij}^*$  (and also update the weights,  $w_{ij}$ , used for estimation in step 1).

This two-step algorithm can be iterated until convergence has been achieved. The resulting estimates are known as the PQL estimates. Recall that in the linear mixed effects model the estimator of  $\beta$  is the *generalized least squares* (GLS) estimator that can be expressed as

$$(15.4) \quad \hat{\beta} = \left\{ \sum_{i=1}^N (X_i' V_i^{*-1} X_i) \right\}^{-1} \sum_{i=1}^N (X_i' V_i^{*-1} Y_i^*),$$

where  $Y_i^*$  is the  $n_i \times 1$  vector of “working” responses (with components  $Y_{ij}^*$ ),  $X_i$  is the corresponding matrix of covariates,  $V_i^* = Z_i G Z_i' + W_i^{-1}$  is the marginal covariance of  $Y_i^*$ ,  $Z_i$  is the matrix of covariates for the random effects, and  $W_i$  is a diagonal weight matrix (with components  $w_{ij}$  along the diagonal). Estimation of the variances (and covariances) of the random effects can be based on the standard ML or REML estimators for a linear mixed effects model applied to the “pseudo-data”. This method of estimation for the variances (and covariances) of the random effects is often referred to as *pseudo-likelihood* (PL).

There are two important points to recognize about the PQL method. First, PQL is an approximate method. This means that in some settings the approximation is quite accurate and produces valid estimates of the fixed effects (and covariance parameters of the random effects), while in other settings it is not accurate and can yield badly biased estimates. In Section 15.4, we consider the factors that influence the accuracy of the approximation. Second, the PQL method produces an estimate of the fixed effects,  $\beta$ , in the GLMM:

$$g\{E(Y_{ij}|b_i)\} = X_{ij}' \beta + Z_{ij}' b_i,$$

albeit with varying degrees of bias that depend on the accuracy of the approximation. Thus, the PQL method is appropriate when the goal of the analysis is to make subject-specific inferences.

Finally, we note that there can be small differences in the implementation of the PQL algorithm in software packages (e.g., PROC GLIMMIX in SAS and the `glmmPQL` package in R and S-Plus). These differences depend on whether the scale parameter  $\phi$  is regarded as fixed and known or as an additional parameter to be estimated from the data at hand. Although for many distributions in the exponential family  $\phi$  is fixed and known (e.g., binomial and Poisson, where  $\phi = 1$ ), empirically the variability of the response often exceeds that predicted by these distributions. Thus, when using the PQL method to fit GLMMs in existing software packages, it may be necessary to override the default options concerning whether the scale parameter  $\phi$  is constrained (e.g.,  $\phi = 1$ ) or included in the PQL estimation as an additional parameter to be estimated from the data.

## 15.3 MARGINAL QUASI-LIKELIHOOD

There is a second approximation that leads to a method known as marginal quasi-likelihood (MQL) estimation. Similar to PQL, MQL can be motivated via an approximation of the GLMM. With MQL the approximation is based on a Taylor series expansion of (15.2) around estimates  $\hat{\beta}$  and around  $b_i = 0$ . For a canonical link function this results in the following approximation to the model given by (15.2):

$$(15.5) \quad Y_{ij}^{**} \approx X'_{ij}\beta + Z'_{ij}b_i + \epsilon_{ij}^{**},$$

where  $Y_{ij}^{**} = v^{-1}(\hat{\mu}_{ij})(Y_{ij} - \hat{\mu}_{ij}) + X'_{ij}\hat{\beta}$  and  $\epsilon_{ij}^{**} = v^{-1}(\hat{\mu}_{ij})\epsilon_{ij}$ . Thus we now also have a standard linear mixed effects model for the transformed response,  $Y_{ij}^{**}$ , with fixed effects  $\beta$ , random effects  $b_i$ , and within-subject errors  $\epsilon_{ij}^{**}$ . Note that the transformed response in MQL,  $Y_{ij}^{**}$ , depends on  $\hat{\mu}_{ij} = g^{-1}(X'_{ij}\hat{\beta})$ , and not on  $\hat{\mu}_{ij}^b = g^{-1}(X'_{ij}\hat{\beta} + Z'_{ij}\hat{b})$ . That is, in contrast to PQL,  $\mu_{ij}$  depends on the *marginal* linear predictor ( $X'_{ij}\beta$ ) instead of the *conditional* linear predictor ( $X'_{ij}\beta + Z'_{ij}b_i$ ). As we will see, this has some important implications for interpretation of the resulting MQL estimates of  $\beta$ .

Recognizing that the “working” response  $Y_{ij}^{**}$  follows a linear mixed effects model, with fixed effects  $\beta$ , and with random effects  $b_i$  and within-subject errors  $\epsilon_{ij}^{**}$ , estimation can proceed, similar to PQL, by iterating between the following two steps:

1. Fit a linear mixed effects model to the working response  $Y_{ij}^{**}$ , to obtain updated estimates of  $\beta$  and  $G$ . The linear mixed effects model is fitted using weights,  $w_{ij}$ , that are inversely proportional to the variance of the  $\epsilon_{ij}^{**}$ ; for example,  $w_{ij} = \phi^{-1} v(\hat{\mu}_{ij})$ , when a canonical link function is adopted.
2. Use the estimates of  $\beta$  obtained in step 1 to update the working response  $Y_{ij}^{**}$  (and also update the weights,  $w_{ij}$ , used for estimation in step 1).

This two-step algorithm can be iterated until convergence has been achieved. The resulting estimates are known as the MQL estimates.

One crucial aspect of the MQL method is often overlooked. Because the MQL method involves an approximation based on a Taylor series expansion of (15.2) around estimates  $\hat{\beta}$  and around  $b_i = 0$ , it does not yield an estimate of  $\beta$  in the GLMM:

$$g\{E(Y_{ij}|b_i)\} = X'_{ij}\beta + Z'_{ij}b_i.$$

Instead, because the transformed response in MQL,  $Y_{ij}^{**}$ , depends on the *marginal* mean  $\hat{\mu}_{ij} = g^{-1}(X'_{ij}\hat{\beta})$ , and not on the *conditional* mean  $\hat{\mu}_{ij}^b = g^{-1}(X'_{ij}\hat{\beta} + Z'_{ij}\hat{b})$ , the MQL algorithm yields an estimate of the regression parameters  $\beta$  in the following marginal model:

$$g\{E(Y_{ij})\} = X'_{ij}\beta.$$

Thus the MQL method is appropriate only when the goal of the analysis is to make population-averaged, not subject-specific, inferences. Some of the statistical literature on GLMMs, and much of the documentation for commercially available software packages for fitting GLMMs, has been less than transparent about the different targets of inference associated with the PQL and MQL methods. In the next section we highlight how the PQL and MQL methods differ in terms of their respective targets of inference, with the PQL method being appropriate when the goal is to make subject-specific inferences whereas the MQL method is appropriate when the goal is to make population-averaged inferences.

## 15.4 CAUTIONARY REMARKS ON THE USE OF PQL AND MQL

We note that PQL and MQL both involve approximations that result in a linearization of the GLMM. Consequently both methods can be implemented through repeated use of statistical software for fitting standard linear mixed effects models. This feature of both methods makes them computationally simple to apply. However, the two methods differ in terms of the locus of the linearization of (15.2): in PQL the expansion of (15.2) is around current estimates  $\hat{\beta}$  and  $\hat{b}_i$ . In a certain sense, this can be considered a “subject-specific” expansion of (15.2). In contrast, in MQL the expansion of (15.2) is around current estimates  $\hat{\beta}$  and around  $b_i = 0$ , the mean of the random effects. In a similar way the latter can be considered a “population-averaged” expansion of (15.2). As we discuss in greater detail below, this has ramifications for the interpretation of  $\beta$  when estimated using the PQL and MQL algorithms.

Next we consider the factors that influence the accuracy of the approximation. In the case of PQL, the accuracy of the approximation improves when the “sufficient statistics” for  $\beta$  and  $b_i$  have approximate normal distributions. The “sufficient statistics” for  $\beta$  and  $b_i$  can be loosely defined as those quantities that provide all of the information in the sample that is useful for estimating  $\beta$  and  $b_i$ . It can be shown that estimation of  $\beta$  relies on the following sufficient statistics:  $\sum_{j=1}^{n_i} X_{ij} Y_{ij}$ . In contrast, estimation of  $b_i$  relies on  $\sum_{j=1}^{n_i} Z_{ij} Y_{ij}$ . Both of these statistics are weighted averages of the  $Y_{ij}$ , where averaging is over the repeated occasions. Thus, in cases where  $Y_{ij}$  is quantitative (e.g., counts from a Poisson distribution with mean  $> 5$ ), we can expect these weighted averages to have approximate normal distributions even when the number of repeated measurements,  $n_i$ , is relatively small. Conversely, when  $Y_{ij}$  is binomial but based on a small denominator, and especially when  $Y_{ij}$  is binary (with binomial denominator of 1), these weighted averages do not have approximate normal distributions unless the number of repeated measurements is very large. This helps explain why PQL often produces very poor estimates of both the fixed effects and the variance components of the random effects when the response variable is binary and there is only a relatively small number of repeated measurements available on each individual. In such settings, PQL can yield seriously biased estimates of effects; in general, the estimates of the fixed effects, but especially the variance components, will be attenuated toward zero. This systematic underestimation of the variance components has been very well documented in the case of binary responses; consequently PQL should be used only in the setting of binomial proportions when the denominator is relatively large and the numerators take on values in the mid-range, for example, binomial proportions where the expected number of successes (or failures) exceeds 5. In principle, as the number of repeated measurements increases, the bias of PQL decreases accordingly.

It is also quite instructive to compare PQL to adaptive Gaussian quadrature. Recall that adaptive Gaussian quadrature can be made as accurate as necessary by increasing the number of quadrature points. Although the relationship between PQL and quadrature methods may not be transparent, it can be shown that a Laplace approximation to the GLMM likelihood corresponds to adaptive Gaussian quadrature with just a single quadrature point. Consequently, given the very close relation between PQL and methods based on a Laplace approximation, it should not be too surprising that this less than optimal choice for the number of quadrature points can lead to seriously biased estimates of effects in GLMMs.

When making statistical inferences about the fixed effects and the variance components, there are important differences between the PQL method and “true” ML estimation based on adaptive Gaussian quadrature. Unlike ML estimation, inferences based on the PQL method cannot rely on the likelihood, e.g., likelihood ratio test statistics are no longer available. Inferences must, however, rely on Wald statistics and confidence intervals. So we must caution the reader that many of the commercially available software packages (e.g., PROC GLIMMIX in SAS) that implement PQL produce output

that includes values for the maximized “log likelihood” and various information criteria (e.g., AIC) based on the maximized “log likelihood”; these cannot be used for inference and should be completely ignored. The reported maximized “log likelihood” applies only to the working response vector  $Y^*_i$ ; it is not the value of the true “log likelihood” for the fitted GLMM to the response vector  $Y_i$ .

In regard to the properties of the MQL method, in much of the statistical literature there is a consensus that MQL produces badly biased estimates of the GLMM parameters unless the variability of the random effects is close to zero. Moreover, unlike PQL, with MQL the bias remains even as the number of repeated measurements increases. What is less clear in the statistical literature, and consequently in documentation for commercially available software packages for fitting GLMMs, is that the MQL method does not yield an estimate of  $\beta$  in the GLMM:

$$g\{E(Y_{ij}|b_i)\} = X'_{ij}\beta + Z'_{ij}b_i.$$

Instead, because MQL is based on a “population-averaged” expansion of (15.2), it produces an estimate of the regression parameters  $\beta$  in the following marginal model:

$$g\{E(Y_{ij})\} = X'_{ij}\beta,$$

that is, a model for the mean response that does not include any random effects but assumes that the marginal mean relationship has the same link function (as for the GLMM). That is, the only similarity in the GLMM and marginal model given above is that the link function,  $g(\cdot)$ , is adopted by both models; otherwise, the “subject-specific” regression parameters in the GLMM and the “population-averaged” regression parameters in the marginal model are discernibly different (see Chapter 16), especially when the variability of the random effects,  $b_i$ , is large. Because the MQL method is based on fitting a marginal model to the longitudinal data, rather than a GLMM, the resulting regression estimates are not expected to be unbiased for the fixed effects in the GLMM. Put simply, when there is a fundamental mismatch between the model of interest and the model that is fit to the data, it cannot be expected that the latter will yield valid estimates of the former. When MQL is used for estimation and inference in the GLMM, the bias with the MQL method arises from the incorrect target of inference (i.e., the *marginal* rather than *conditional* mean) that is modeled. This is because the MQL method estimates the regression parameters in the marginal model,

$$g\{E(Y_{ij})\} = X'_{ij}\beta,$$

and not the fixed effects in the GLMM,

$$g\{E(Y_{ij}|b_i)\} = X'_{ij}\beta + Z'_{ij}b_i.$$

Indeed, it is false to claim MQL to be a legitimate method of estimation for GLMMs! So we caution the reader that many of the leading commercially available software packages for fitting GLMMs include MQL as an optional method of estimation.

Nevertheless, while we cannot recommend the use of MQL for estimation of effects in a GLMM, we note one important potential use of the MQL algorithm: estimation of parameters in a marginal model for the longitudinal response. That is, in settings where it is somewhat more natural or convenient to express the “working covariance” in a marginal model via the introduction of random effects, MQL can be used to estimate the marginal regression parameters. In this sense, MQL can be regarded as a potentially flexible way to model the “working covariance” in the generalized estimating equations (GEE) approach. We illustrate this use of MQL for estimation of marginal model parameters in Section 15.5.

In summary, the PQL and MQL methods of estimation for GLMMs are based on discernibly different expansions of (15.2) and have distinct targets of inference. Because the PQL method is based on a “subject-specific” expansion of (15.2), it yields estimates of  $\beta$  from the following model for the *conditional* mean:

$$E(Y_{ij}|b_i) = g^{-1}(X'_{ij}\beta + Z'_{ij}b_i),$$

where  $g^{-1}(\cdot)$  denotes the inverse link function. The PQL method yields valid estimates of  $\beta$  provided that  $Y_{ij}$  is quantitative (e.g., counts from a Poisson distribution with mean greater than 5 or binomial proportions where the expected number of successes exceeds 5). PQL can yield badly biased estimates of  $\beta$  (and the variance components) when the response variable is binary and there are few

repeated measurements. In the latter settings, we cannot recommend the use of the PQL method. In contrast to PQL, the MQL method is based on a “population-averaged” expansion of (15.2), and it yields estimates of  $\beta$  from the following model for the *marginal* mean (averaged over the distribution of the random effects):

$$E(Y_{ij}) = g^{-1}(X'_{ij}\beta),$$

with  $\text{Var}(Y_{ij}) = v\{E(Y_{ij})\} \phi$ . As discussed in Chapters 12 and 13, the components of  $\beta$  in this marginal model have “population-averaged” interpretations. Because the subject-specific and population-averaged regression parameters can be discernibly different, we cannot recommend that the MQL method be used for estimation of the fixed effects in a GLMM. Instead, we see the MQL method as having potential use for the estimation of regression parameters in marginal models only. Finally, we note that if a marginal model, rather than a GLMM, is specified for the vector of responses, then the PQL and MQL methods are formally equivalent because they are based on the same expansion. In that setting, both methods can be considered GEE estimators of the marginal model parameters.

To help the reader digest the implications of the choice between the PQL and MQL methods for statistical inference, we present a summary in [Table 15.1](#) that contrasts the true targets of inference (i.e., “subject-specific” or “population-averaged”) for the estimated regression parameters versus model specification (GLMM versus marginal model) and method of estimation (PQL versus MQL). From [Table 15.1](#) it is apparent that only the PQL method yields estimates of the fixed effects from a GLMM when a GLMM has been specified for the data at hand. In contrast, when a GLMM has been specified, the MQL method estimates the regression parameters in a marginal model for the vector of responses that has the same link (and variance) function as the conditional distribution of the response in the specified GLMM. Finally, in cases where a marginal model has been specified for the data, the resulting estimates of the marginal regression parameters are valid regardless of the choice of estimation method; indeed, for a marginal model specification, the PQL and MQL methods yield identical estimators of the marginal regression parameters.

**Table 15.1** Targets of inference for the estimated regression parameters as a function of model specification (GLMM versus marginal model) and method of estimation (PQL versus MQL).

Model Specification		
Estimation Method	GLMM	Marginal
PQL	Subject-specific	Population-averaged
MQL	Population-averaged <sup>a</sup>	Population-averaged

<sup>a</sup>Yields estimates of regression parameters from marginal model with the same link and variance functions as specified in the GLMM.

Many of the issues highlighted in this section can be solidified by a numerical illustration based in part on the example given in Section 14.3 (see [Figure 14.1](#)). Assume that  $N$  individuals are measured repeatedly at baseline and after treatment with a new drug intended to reduce the risk of disease. Specifically,  $n/2$  measurements of the response are obtained at baseline, and  $n/2$  measurements post-baseline. To model changes in the response probabilities from baseline to post-baseline, we consider the following logistic regression model, with normally distributed random intercepts:

$$\text{logit}\{E(Y_{ij}|b_i)\} = \beta_1^* + \beta_2^* t_{ij} + b_i,$$

where  $t_{ij} = 0$  at baseline and  $t_{ij} = 1$  post-baseline ( $i = 1, \dots, N$ ;  $j = 1, \dots, n$ ). For a “typical” individual from the population (where a “typical” individual is one with unobserved random effect  $b_i = 0$ ), the log odds of disease at baseline is  $\beta_1^*$ ; the log odds of disease following treatment with the new drug is  $\beta_1^* + \beta_2^*$ .

We can simulate data from this model for different values of  $N$  and  $n$  for the case where  $\beta_1^* = 1.5$ ,  $\beta_2^* = -3.0$ , and  $\text{Var}(b_i) = 1.0$ . Then we estimate the model parameters using PQL and MQL methods. Recall that the PQL method yields estimates of  $\beta_1^*$  and  $\beta_2^*$ , albeit potentially biased estimates. In contrast, the MQL method yields estimates of the regression parameters from the marginal logistic

model,

$$\text{logit}\{E(Y_{ij})\} = \beta_1 + \beta_2 t_{ij}.$$

Because of the discreteness of the only covariate ( $t_{ij} = 0$  or  $1$ ) in the logistic regression model with random intercepts, a logistic regression model also holds for the marginal probabilities (averaged over the distribution of the random effects). However, the regression parameters in the marginal model are attenuated, with  $\beta_1 = 1.257$  and  $\beta_2 = -2.514$ . In [Table 15.2](#) we present the PQL and MQL estimates of the “fixed effects” and variance component  $\text{Var}(b_i)$  (for the latter, PQL estimates are presented only) when  $N = 5000$  and  $n = 4, 8, 16, 32, 64$ , and  $128$ . We have purposely chosen a very large value for  $N$ , so that any concerns about sampling variation of the reported estimates can be completely set aside, but have allowed for a broad range of values for  $n$ , the number of repeated measures.

**Table 15.2** PQL and MQL estimates, and percent relative bias (in parentheses) when compared to true parameter values, as a function of the number of repeated measurements ( $n$ ) for simulated data from the logistic regression model with randomly varying intercepts.

True Value	$\beta_1^*$	$\beta_1$	$\beta_2^*$	$\beta_2$	$\text{Var}(b_i)$
	PQL	MQL	PQL	MQL	PQL
4	1.310	1.257	-2.619	-2.513	0.647
	(12.7%)	(0.0%)	(12.7%)	(0.0%)	(35.3%)
8	1.370	1.267	-2.752	-2.545	0.780
	(8.7%)	(0.8%)	(8.3%)	(1.2%)	(22.0%)
16	1.415	1.263	-2.844	-2.539	0.864
	(5.7%)	(0.5%)	(5.2%)	(1.0%)	(13.6%)
32	1.443	1.255	-2.902	-2.511	0.947
	(3.8%)	(0.6%)	(3.3%)	(0.1%)	(5.3%)
64	1.470	1.255	-2.939	-2.510	0.962
	(2.0%)	(0.2%)	(2.0%)	(0.1%)	(3.8%)
128	1.490	1.263	-2.972	-2.519	0.971
	(0.7%)	(0.5%)	(0.9%)	(0.2%)	(2.9%)

Note:  $\beta_1^*$  and  $\beta_2^*$  are the regression parameters in the random effects logistic regression model;  $\beta_1$  and  $\beta_2$  are the corresponding regression parameters in the marginal logistic regression model (averaged over the distribution of the random effects).

[Table 15.2](#) shows how the PQL and MQL methods yield estimates of the regression parameters from GLMM and marginal models, respectively. In the last row of [Table 15.2](#) we can see that the PQL and MQL estimates are very close to the true values of the parameters in the respective random effects and marginal models for these data. This confirms our earlier warning about the different targets of inference for these two different methods of estimation. [Table 15.2](#) also shows that the MQL method yields unbiased estimates of the marginal regression parameters. The relative bias of the MQL estimates of the marginal regression parameters is negligible and less than 1% in most instances. In contrast, the PQL estimates of the fixed effects and variance component are badly biased when the number of repeated measures ( $n$ ) is relatively small. For example, when  $n = 4$ , the relative bias of the estimates of the fixed effects is approximately 13%, while the relative bias of the estimate of  $\text{Var}(b_i)$  is approximately 35%. The bias of the estimates of the fixed effects remains, and exceeds 5%, even for  $n = 16$ . Recall that the marginal probability of success (here denoting disease) at baseline is approximately 0.78. Thus for  $n = 16$  the binomial denominator at baseline is  $n/2 = 8$ , and the expected number of successes and/or failures ( $0.78 \times 8 = 6.2$  or  $0.22 \times 8 = 1.8$ ) does not exceed 5. So, in light of our earlier warnings about the use of PQL with binary data, it should not be too surprising that there is discernible bias when  $n = 16$  (and the expected number of failures is less than 2). In contrast, when  $n = 64$  the binomial denominator at baseline is  $n/2 = 32$ , and the expected number of failures ( $0.22 \times 32 = 7.04$ ) now exceeds 5; for  $n = 64$  the bias is almost negligible (approximately 2% for the estimates of the fixed effects). The pattern of results for the PQL method highlights how the bias diminishes as the number of repeated measurements increases.

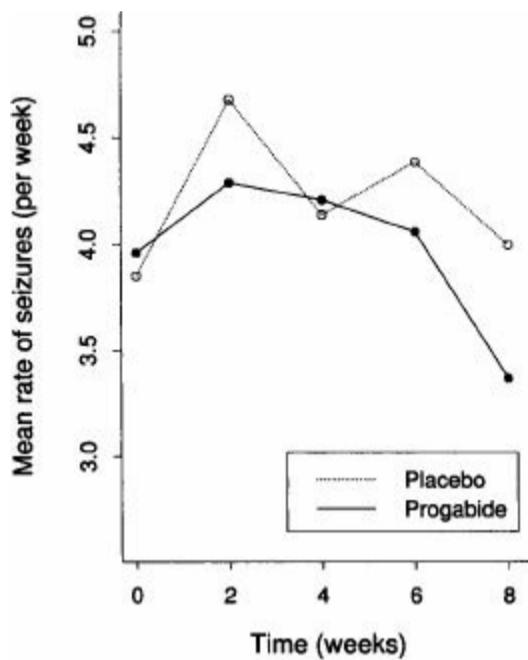
## 15.5 CASE STUDIES

In this section we illustrate the main ideas presented in this chapter by considering generalized linear mixed effects models for analyzing longitudinal data from two studies. The first illustration considers a Poisson regression model, with random effects, for analyzing count data on epileptic seizures from a clinical trial of the anti-epileptic drug, pro gabide. The second illustration considers a logistic regression model, with random effects, for analyzing data on amenorrhea from a randomized clinical trial of contracepting women. These two data sets were previously analyzed in Section 14.7 using GLMMs fit via adaptive Gaussian quadrature. Here we fit similar models for the conditional means using the PQL method to highlight settings where it should and should not be applied. We also present analyses using the MQL method to highlight how this method yields estimates of regression parameters from a marginal rather than conditional (or subject-specific) model.

# Clinical Trial of an Anti-epileptic Drug

The first example involves data from the placebo-controlled clinical trial of 59 epileptic patients, conducted by Leppik et al. (1987). Patients with partial seizures were randomized to either progabide or a placebo, as an adjuvant to the standard anti-epileptic therapy. Prior to treatment the number of epileptic seizures during the preceding 8-week interval was recorded. Counts of epileptic seizures during 2-week intervals before each of four successive post-randomization clinic visits were also recorded. The average rates of seizures (per week) at baseline and in the four post-randomization visits are displayed in [Figure 15.1](#). Note that the mean rates over time in the two treatment groups vary from approximately 3 seizures per week to approximately 5 seizures per week. Thus the expected count of the number of seizures in any 2-week interval is approximately 6 to 10. As we will discuss later, the fact that the mean number of seizures is relatively large has implications for the accuracy of the approximation used in the PQL method.

**Fig. 15.1** Mean rate of seizures (per week) at baseline, week 2, week 4, week 6, and week 8 in the progabide and placebo groups.



As was noted in Section 14.7, patient 49 is an outlier with extreme counts at all occasions. For illustrative purposes we exclude observations on this patient ( $N=58$ ) and replicate the analysis displayed in [Table 14.7](#). That is, we consider the following mixed effects log-linear regression model for the subject-specific expected counts (or rates) of seizures:

$$\begin{aligned} \log E(Y_{ij}|b_i) &= \log(T_{ij}) + (\beta_1 + b_{1i}) + (\beta_2 + b_{2i}) \text{Time}_{ij} + \beta_3 \text{Trt}_i \\ &\quad + \beta_4 \text{Trt}_i \times \text{Time}_{ij}, \end{aligned}$$

where  $Y_{ij}$  is the number of epileptic seizures for the  $i^{th}$  patient in the  $j^{th}$  period of observation ( $j = 0, \dots, 4$ ), and  $T_{ij}$  is the length of period  $j$  (where  $T_{ij} = 8$  if  $j = 0$  and  $T_{ij} = 2$  if  $j = 1, 2, 3, 4$ ). The offset,  $\log(T_{ij})$ , is included because the “time at risk” is not the same in the baseline (8 weeks) and four successive follow-up periods (2-week intervals). The variable Trt is an indicator variable for treatment group, with Trt = 0 if an individual was randomized to the placebo group and Trt = 1 if randomized to the progabide group. The binary variable Time denotes the baseline and follow-up periods, with Time = 0 for the baseline period and Time = 1 for the follow-up periods (periods 1–4). Given  $b_i$ , it is assumed that the  $Y_{ij}$  are independent and have a Poisson distribution, with  $\text{Var}(Y_{ij}|b_i) = E(Y_{ij}|b_i)$ . Initially, we assume Poisson variation for  $\text{Var}(Y_{ij}|b_i)$  and fix  $\phi = 1$ . In subsequent analyses, we include  $\phi$  as a parameter to be estimated from the data, allowing for potential overdispersion. Finally, we assume that the random intercepts and slopes,  $b_i$ , have a bivariate normal distribution, with zero mean and  $2 \times 2$  covariance matrix  $G$ .

The PQL estimates of the fixed effects and covariance parameters for the log-linear model are displayed in [Table 15.3](#). A test of the null hypothesis,  $H_0: \beta_4 = 0$ , indicates that there is a significant time  $\times$  treatment interaction at the 0.05 level. These results suggest that there are differences between

the two treatments in terms of subject-specific changes in the expected rates of seizures, with a greater reduction in the expected seizure rate from baseline for a typical patient treated with progabide. Specifically, for a patient receiving the placebo, there is almost no expected change in the rate of seizures (or  $1 - e^{0.0115} \approx 0$ ), while for a patient treated with progabide the expected decrease in seizures is approximately 30% (or  $1 - e^{0.0115-0.3415} \approx 0.28$ ).

**Table 15.3** PQL estimates and standard errors from mixed effects log-linear regression model for the seizure data, excluding patient 49.

Variable	Estimate	SE	Z
Intercept	1.0869	0.1344	8.09
Time	0.0115	0.1058	0.11
Trt	-0.0074	0.1868	-0.04
Trt × Time	-0.3415	0.1490	-2.29
$g_{11} = \text{Var}(b_{1i})$	0.4579	0.0954	
$g_{22} = \text{Var}(b_{2i})$	0.2196	0.0594	
$g_{12} = \text{Cov}(b_{1i}, b_{2i})$	0.0122	0.0537	

Note: PQL estimation assuming  $\phi$  fixed at 1.

It is instructive to compare the PQL estimates of the fixed effects and covariance parameters to the corresponding “true” ML estimates (obtained using 50-point adaptive Gaussian quadrature) in [Table 14.7](#). In general, there is close agreement, indicating that the PQL method provides an adequate approximation in this setting. This is to be expected because, as was noted earlier, the mean of the seizure counts at each occasion is relatively large. Thus for this particular data set the conditions required for the appropriate use of the PQL method are met. We also repeated the PQL analysis but included  $\phi$  as a parameter to be estimated from the data, allowing for potential overdispersion. This yielded an estimate of overdispersion,  $\hat{\phi} = 1.94$ , suggesting that the conditional variability in the data may be twice as large as that predicted by Poisson variability. However, it should be mentioned that the estimates of the variances of the random effects were noticeably smaller than in [Table 15.3](#), perhaps reflecting that  $\phi$  and the variance component parameters are in competition with each other to explain the variability in the data. The estimate of the fixed effect for the time  $\times$  treatment interaction,  $\hat{\beta}_4 = -0.303$  (SE = 0.151), is comparable in magnitude to the estimate in [Table 15.3](#) and leads to the same conclusion about the benefits of progabide.

Finally, we re-fit the original model (assuming  $\phi = 1$ ) except now using the MQL method. The results of this analysis are presented in [Table 15.4](#). Note that the estimates of the fixed effects are discernibly different from those displayed in [Table 15.3](#). For example, the estimate of the effect of greatest interest,  $\beta_4$ , is approximately 12% smaller in absolute value. Conversely, the estimate of the intercept is approximately 35% larger. These differences in the estimates of the fixed effects are not a reflection of sampling variability; rather, they reflect the different targets of inference of the PQL and MQL methods. That is, the MQL method does not yield estimates of the fixed effects in the model,

$$\begin{aligned}\log E(Y_{ij}|b_i) &= \log(T_{ij}) + (\beta_1 + b_{1i}) + (\beta_2 + b_{2i}) \text{Time}_{ij} + \beta_3 \text{Trt}_i \\ &\quad + \beta_4 \text{Trt}_i \times \text{Time}_{ij};\end{aligned}$$

**Table 15.4** MQL estimates and standard errors from mixed effects log-linear regression model for the seizure data, excluding patient 49.

Variable	Estimate	SE	Z
Intercept	1.3476	0.1574	8.56
Time	0.1118	0.1159	0.96
Trt	-0.1068	0.1937	-0.55
Trt × Time	-0.3024	0.1711	-1.77
$g_{11} = \text{Var}(b_{1i})$	0.5182	0.1043	
$g_{22} = \text{Var}(b_{2i})$	0.3697	0.0834	
$g_{12} = \text{Cov}(b_{1i}, b_{2i})$	-0.0127	0.0660	

Note: MQL estimation assuming  $\phi$  fixed at 1; SE based on empirical variance estimator.

instead, it yields estimates of the regression parameters in the marginal model,

$$\log E(Y_{ij}) = \log(T_{ij}) + \beta_1 + \beta_2 \text{Time}_{ij} + \beta_3 \text{Trt}_i + \beta_4 \text{Trt}_i \times \text{Time}_{ij},$$

that has the same log link function and variance function,  $\text{Var}(Y_{ij}) = E(Y_{ij})$ , and with a “working” covariance determined via the introduction of randomly varying intercepts and slopes. To reinforce this point, we used the standard GEE method to fit the marginal model given above to the epilepsy data. The GEE estimates of the regression parameters, obtained under an unstructured “working” covariance, are presented in [Table 15.5](#). As expected, the estimated regression parameters are very close to those obtained using the MQL method. The very small differences in the estimates of the regression parameters are due to the different choices of “working” covariance. In general, with complete data and relatively large sample sizes, we would expect to find no important differences between the estimates from the MQL method and the more standard GEE approach. Indeed, when the MQL method is applied to a GLMM, it should simply be regarded as a GEE method of estimation of the parameters in a marginal model that paradoxically specifies the “working” covariance through the introduction of random effects. Thus, when applying the MQL method to a GLMM, the specification of random effects only has implications for the implicit choice of a “working” covariance for the corresponding marginal model that assumes the same link and variance functions for the marginal means and variances, respectively.

[Table 15.5](#) GEE estimates and standard errors from marginal log-linear regression model for the seizure data, excluding patient 49.

Variable	Estimate	SE	Z
Intercept	1.3312	0.1611	8.26
Time	0.1141	0.0926	1.23
Trt	-0.1021	0.1959	-0.52
Trt × Time	-0.3167	0.1494	-2.12

Note: GEE estimation with unstructured “working” covariance; SE based on empirical variance estimator.

# Clinical Trial of Contracepting Women

The next example is from a longitudinal clinical trial of contracepting women reported by Machin et al. (1988). In this trial women received an injection of either 100 mg or 150 mg of depot-medroxyprogesterone acetate (DMPA) on the day of randomization and three additional injections at 90-day intervals. There was a final follow-up visit 90 days after the fourth injection. Throughout the study each woman completed a menstrual diary that recorded any vaginal bleeding pattern disturbances. The diary data were used to determine whether a woman experienced amenorrhea, the absence of menstrual bleeding, for a specified number of days.

A total of 1151 women completed the menstrual diaries, and the diary data were used to generate a binary sequence for each woman, according to whether she had experienced amenorrhea in the four successive three-month intervals. A feature of this clinical trial is that there was substantial dropout. More than one-third of the women dropped out before the completion of the trial; 17% dropped out after receiving only one injection of DMPA, 13% dropped out after receiving only two injections, and 7% dropped out after receiving three injections. The outcome of interest is a binary response indicating whether a woman experienced amenorrhea in the four successive three-month intervals. The goal of the analyses presented here is to determine subject-specific changes in the risk of amenorrhea over the course of the study (12 months), and the influence of dosage of DMPA on changes in a woman's risk of amenorrhea.

Letting  $Y_{ij} = 1$  if the  $i^{th}$  woman experienced amenorrhea in the  $j^{th}$  injection interval ( $j = 1, \dots, 4$ ), and  $Y_{ij} = 0$  otherwise, we consider the following mixed effects logistic regression model for  $Y_{ij}$ ,

$$\begin{aligned}\text{logit}\{E(Y_{ij}|b_i)\} &= \beta_1 + \beta_2 \text{Time}_{ij} + \beta_3 \text{Time}_{ij}^2 + \beta_4 \text{Dose}_i \times \text{Time}_{ij} \\ &\quad + \beta_5 \text{Dose}_i \times \text{Time}_{ij}^2 + b_i,\end{aligned}$$

where  $\text{Time} = 1, 2, 3, 4$  for the four consecutive 90-day injection intervals, and  $\text{Dose} = 1$  if randomized to 150 mg of DMPA, and  $\text{Dose} = 0$  otherwise. Due to randomization, we assume that the baseline risk (at  $\text{Time} = 0$ ) is the same in both dosage groups and omit a main effect of dose from the model. Given  $b_i$ , it is assumed that the  $Y_{ij}$  are independent and have a Bernoulli distribution, with  $\text{Var}(Y_{ij}|b_i) = E(Y_{ij}|b_i) \{1 - E(Y_{ij}|b_i)\}$ , and  $\phi = 1$ . Finally, we assume that the single random effect  $b_i$  has a univariate normal distribution, with zero mean and variance  $\sigma_b^2$ ,  $b_i \sim N(0, \sigma_b^2)$ .

The PQL estimates of the regression parameters for this model are presented in [Table 15.6](#). These results provide evidence that the subject-specific log odds of amenorrhea increases over the 12 months of the trial, and that subject-specific changes in the risk of amenorrhea depend on the dose of DMPA. For example, for a woman assigned to the low dose of DMPA, the log odds of amenorrhea increases approximately linearly, with an increase in the log odds of 0.75 (or  $0.7735 - 0.0267$ ) at 3 months, 1.44 (or  $2 \times 0.7735 - 4 \times 0.0267$ ) at 6 months, 2.08 (or  $3 \times 0.7735 - 9 \times 0.0267$ ) at 9 months, and 2.67 (or  $4 \times 0.7735 - 16 \times 0.0267$ ) at 12 months. These increases in risk correspond to odds ratios of 2.1 (or  $e^{0.75}$ ), 4.2 (or  $e^{1.44}$ ), 8.0 (or  $e^{2.08}$ ), and 14.4 (or  $e^{2.67}$ ) at 3, 6, 9, and 12 months, respectively. Note that the estimates of the fixed effects are discernibly smaller in absolute value than the “true” ML estimates reported in [Table 14.2](#). The difference between the estimates is even more marked for the variance of the random effect. The ML estimate is approximately 5.1 whereas the PQL estimate is approximately 1.8; the former is almost 3 times larger than the latter. This highlights how badly biased the PQL estimates can be when the response is binary and there are relatively few repeated measures. In this setting the PQL method cannot be recommended and inferences based on the PQL estimates are not trustworthy.

**Table 15.6** PQL estimates and standard errors from a mixed effects logistic regression model, with randomly varying intercepts, for the amenorrhea data.

Variable	Estimate	SE	Z
Intercept	-2.5859	0.2298	-11.25
Time	0.7735	0.2178	3.55
Time <sup>2</sup>	-0.0267	0.0449	-0.59
Dose × Time	0.3666	0.1416	2.59
Dose × Time <sup>2</sup>	-0.0725	0.0385	-1.88
$\sigma_b^2$	1.8488	0.1693	

Note: PQL estimation assuming  $\phi$  is fixed at 1.

Finally, for illustrative purposes only, we also fit the same model using the MQL method. The results of this analysis are presented in [Table 15.7](#). In general, the estimates of the fixed effects and the variance component are attenuated toward zero (in absolute value) relative to the corresponding estimates produced by the PQL method. However, we remind the reader that the PQL method yields estimates, albeit biased estimates, of the regression parameters from the model for the conditional log odds,

$$\begin{aligned}\text{logit}\{E(Y_{ij}|b_i)\} = & \beta_1 + \beta_2 \text{Time}_{ij} + \beta_3 \text{Time}_{ij}^2 + \beta_4 \text{Dose}_i \times \text{Time}_{ij} \\ & + \beta_5 \text{Dose}_i \times \text{Time}_{ij}^2 + b_i.\end{aligned}$$

**Table 15.7** MQL estimates and standard errors from a mixed effects logistic regression model, with randomly varying intercepts, for the amenorrhea data.

Variable	Estimate	SE	Z
Intercept	-2.2163	0.1770	-12.52
Time	0.6851	0.1598	4.29
Time <sup>2</sup>	-0.0314	0.0322	-0.97
Dose × Time	0.3080	0.1115	2.76
Dose × Time <sup>2</sup>	-0.0611	0.0290	-2.11
$\sigma_b^2$	1.2608	0.1184	

Note: MQL estimation assuming  $\phi$  is fixed at 1; SE based on empirical variance estimator.

In contrast, the MQL method yields perfectly valid estimates of the regression parameters from the following model for the marginal log odds:

$$\begin{aligned}\text{logit}\{E(Y_{ij})\} = & \beta_1 + \beta_2 \text{Time}_{ij} + \beta_3 \text{Time}_{ij}^2 + \beta_4 \text{Dose}_i \times \text{Time}_{ij} \\ & + \beta_5 \text{Dose}_i \times \text{Time}_{ij}^2.\end{aligned}$$

# 15.6 COMPUTING: FITTING GENERALIZED LINEAR MIXED MODELS USING PROC GLIMMIX IN SAS

PROC GLIMMIX in SAS is a procedure for fitting GLMMs (among other models) using maximum likelihood (via adaptive quadrature), PQL, or MQL methods; in this section we focus on PQL and MQL methods only. The procedure is very versatile and the syntax is remarkably similar to that used in PROC MIXED. For example, PROC GLIMMIX has a RANDOM statement that is used in a way similar to PROC MIXED for introducing random effects. Although PROC GLIMMIX does not have an explicit REPEATED statement, it does include an option on the RANDOM statement that mirrors how the REPEATED statement is used in PROC MIXED. Recall that the algorithms for both the PQL and MQL methods require repeated use of software for fitting linear mixed effects models, so this close correspondence in syntax should not be too surprising.

Of note, much of the initial output produced by PROC GLIMMIX is standard output from a linear mixed effects model applied to the working response vector,  $\mathbf{Y}_i^*$  or  $\mathbf{Y}_i^{**}$ . Inferences should be based only of the reported estimates, standard errors, and Wald statistics. In particular, the reader is cautioned that the reported value of the (pseudo) log-likelihood, and various (pseudo) likelihood-based goodness of fit statistics, in the PQL or MQL output should not be used for making inferences.

Before discussing syntax in any detail, it is important to note that PROC GLIMMIX allows the specification of two distinct types of models and two distinct methods of estimation of the model parameters. Specifically, we can distinguish models in terms of the inclusion or absence of random effects. The former are referred to as “conditional” (i.e., conditional on random effects) or “subject-specific” models; the latter are referred to as “marginal” or “population-averaged” models. Thus PROC GLIMMIX can be used for fitting both GLMMs and marginal models to longitudinal data. PROC GLIMMIX provides two approximate methods of estimation of the model parameters: PQL and MQL. It also provides a method based on the Laplace approximation; as noted earlier in the chapter, a Laplace approximation to the integrand in the likelihood for a GLMM corresponds to the use of adaptive Gaussian quadrature with only a single quadrature point.

If a marginal model has been specified, both PQL and MQL can be considered GEE approaches, and the two methods will yield identical estimates of the marginal regression parameters and the “working” covariance. So, in cases where a marginal model has been specified, there is no distinction between the PQL and MQL methods. However, in cases where a GLMM has been specified, there are important differences between PQL and MQL. As discussed in earlier sections of this chapter, only the PQL method yields estimates, albeit sometimes biased estimates, of the fixed effects in a GLMM. The MQL method emphatically does not yield estimates of the fixed effects in a GLMM; instead, it provides estimates of the regression parameters in a marginal model that has the same assumed link and variance functions (e.g., logit link function and Bernoulli variance for the analysis of binary responses) for the marginal means and variances, respectively.

In distinguishing PQL and MQL we note also that PROC GLIMMIX provides two small variations of these two methods of estimation: “restricted maximum” and “maximum” PQL and MQL estimation. As noted in Sections 15.2 and 15.3, the PQL and MQL algorithms solve the linear mixed effects model likelihood equations for a “working” vector of responses, denoted by  $\mathbf{Y}_i^*$  or  $\mathbf{Y}_i^{**}$ , respectively; thus, as with standard linear mixed effects models, estimation of the covariance of the random effects (or the marginal covariance) can be based on either the REML or ML equations for these parameters. “Restricted maximum” PQL and MQL refer to the algorithms that use the REML equations for estimation of the covariance parameters; “maximum” PQL and MQL refer to the algorithms that use the ML equations. In PROC GLIMMIX, the former are denoted RSPL and RMPL, the latter are denoted MSPL and MMPL. These abbreviations are not very intuitive but can be deciphered as follows. The last two letters, PL, denote that the PQL and MQL methods of estimation are based on a

so-called pseudo-likelihood for a linearization (or approximation) to the model. The first letter distinguishes between “restricted maximum” (R) and “maximum” (M) PQL and MQL estimation of the covariance parameters. Finally, the second letter distinguishes between the PQL (S) and MQL (M) methods, recognizing that the former method is based on a “subject-specific” (S) expansion of the model, while the latter is based on a “marginal” (M) or “population-averaged” expansion.

The possible choices of methods of estimation, and their implications for inference (e.g., “subject-specific” or “population-averaged”), are summarized in [Table 15.8](#). For example, the fixed effects (and variance components) in a GLMM can be estimated using either restricted or maximum PQL methods through the use of the METHOD=RSPL or METHOD=MSPL options, respectively. In the absence of random effects in the model, there is no distinction between PQL and MQL methods for estimating the parameters in a marginal model. Thus, when PROC GLIMMIX is used for fitting marginal models, the choice of either METHOD=RSPL or METHOD=RMPL yields identical estimates and standard errors; similarly the choice of either METHOD=MSPL or METHOD=MMPL produces identical results for fitting marginal models.

**Table 15.8** Targets of inference for the estimated regression parameters as a function of model specification (GLMM versus marginal model) and method of estimation (PQL versus MQL) using the METHOD=<option> in PROC GLIMMIX in SAS.

Estimation Method	Model Specification	
	GLMM	Marginal
PQL	Subject-specific	Population-averaged
Restricted Maximum PQL	(RSPL)	(RSPL or RMPL)
Maximum PQL	(MSPL)	(MSPL or MMPL)
MQL	Population-averaged <sup>a</sup>	Population-averaged
Restricted Maximum MQL	(RMPL)	(RMPL or RSPL)
Maximum MQL	(MMPL)	(MMPL or MSPL)

<sup>a</sup>Yields estimates of regression parameters from a marginal model with the same link and variance functions as specified in the GLMM.

One appealing feature of PROC GLIMMIX is that it can fit a broad and flexible class of models. To do so, PROC GLIMMIX makes a distinction between two sources of variation and covariation in the model for the data: (1) variation due to random effects,  $b_i$ , and (2) “residual” (co)variation. To distinguish these two sources of (co)variation, PROC GLIMMIX refers to (1) as “G-side” effects and (2) as “R-side” effects. These two non-standard terms are derived from notation for the covariance matrices for the random effects (denoted G) and the “residual errors” (denoted R) in PROC MIXED for linear mixed effects models. PROC GLIMMIX is quite versatile in allowing models to be fit with various combinations of “G-side” and “R-side” effects. For example, in a standard marginal model, there are no random effects  $b_i$ ; consequently “G-side” effects are completely absent, and the marginal covariance is ordinarily specified in terms of “R-side” effects only. In a standard GLMM, the introduction of random effects is handled by including “G-side” effects, and the “conditional independence” assumption is handled by assuming a simple structure for the “R-side” effects (i.e., uncorrelated residual errors, the default option for the “R-side” effects). In principle, it is possible to fit models that relax the “conditional independence” assumption by allowing for a more general structure for the “R-side” effects (e.g., autoregressive residual errors); however, we caution that, as with linear mixed effects models, there can be subtle issues of model identification when a more general structure for the “R-side” effects is assumed because it may not be possible to estimate both the “G-side” and the more general “R-side” effects from the data at hand.

Because the “G-side” and “R-side” terminology has been borrowed from PROC MIXED, it is useful to compare how PROC MIXED and PROC GLIMMIX specify these effects. In PROC MIXED, the RANDOM statement is used to specify the random effects  $b_i$  (the “G-side” effects). In PROC MIXED, multiple RANDOM statements are allowed. Similarly, in PROC GLIMMIX, the

RANDOM statement is used to specify the random effects  $b_i$  (“G-side” effects) in a GLMM; PROC GLIMMIX also allows the use of multiple RANDOM statements. In PROC MIXED, the REPEATED statement is used to specify assumptions about the “residual errors” (“R-side” effects) in models with or without random effects. This is where the syntax for PROC GLIMMIX departs from PROC MIXED. PROC GLIMMIX does not have a REPEATED statement. Instead, it includes a refinement to the RANDOM statement that makes it completely equivalent to a REPEATED statement, via the use of the RESIDUAL option or through the inclusion of a \_RESIDUAL\_ effect (these options will be explained in detail below). So, unlike in PROC MIXED where both a RANDOM and REPEATED statement can be used in specifying a linear mixed effects model, in PROC GLIMMIX assumptions about the random effects and the “residual errors” are made using two RANDOM statements, with the second being a variant of the RANDOM statement that includes an option or reserved keyword that signifies that it is referring to the “R-side” effects rather than the “G-side” effects.

In many other respects, though, PROC MIXED and PROC GLIMMIX have a lot of similarity in terms of command syntax. For example, to fit a logistic regression model with randomly varying intercepts to longitudinal data from two groups using “restricted” PQL, we can use the illustrative SAS commands given in [Table 15.9](#). Similarly, to fit a mixed effects log-linear regression, with randomly varying intercepts and slopes (also via “restricted” PQL), we can use the illustrative SAS commands given in [Table 15.10](#).

**Table 15.9** Illustrative commands for a mixed effects logistic regression, with randomly varying intercepts, fitted using PQL in PROC GLIMMIX in SAS.

---

```
PROC GLIMMIX METHOD=RSPL;
  CLASS id group;
  MODEL y=group time group*time / DIST=BINOMIAL LINK=LOGIT S;
  RANDOM INTERCEPT / SUBJECT=id TYPE=UN;
```

---

**Table 15.10** Illustrative commands for a mixed effects log-linear regression, with randomly varying intercepts and slopes, fitted using PQL in PROC GLIMMIX in SAS.

---

```
PROC GLIMMIX METHOD=RSPL;
  CLASS id group;
  MODEL y=group time group*time / DIST=POISSON LINK=LOG S;
  RANDOM INTERCEPT time / SUBJECT=id TYPE=UN;
```

---

To fit a marginal logistic regression model, rather than a mixed effects model, to longitudinal data from two groups, we can use the illustrative SAS commands given in [Table 15.11](#); note that for a marginal model specification, the options METHOD=RSPL and METHOD=RMPL yield identical results. This is the most direct and transparent way of fitting a marginal model, since the absence of any random effects is denoted by the use of the RESIDUAL option on the single RANDOM statement. However, recognizing that the MQL method is simply another GEE method for estimating marginal model parameters, we can also fit the marginal logistic regression model using the illustrative SAS commands given in [Table 15.12](#). Note that although the model is specified to be a GLMM with randomly varying intercepts, MQL estimates are requested by using the option METHOD=RMPL instead of METHOD=RSPL. This implies that the resulting MQL estimates of the regression parameters refer not to the logistic regression model with randomly varying intercepts but to the marginal logistic regression model. Although this approach to estimation of marginal model parameters is somewhat less transparent, because it invokes a GLMM, it may have potential applications in settings where the data are highly unbalanced over time. With inherently unbalanced data, it is somewhat more appealing to specify a “working” covariance via the inclusion of randomly varying intercepts and slopes; that is, the “working” covariance is expressed as an explicit function of the times of measurement when times of measurement are included in the covariates for the random effects ( $Z_i$ ). Next we present a brief description of each of the command statements used in [Tables 15.9](#) through [15.12](#).

**Table 15.11** Illustrative commands for a marginal logistic regression model, with exchangeable working correlation, fitted using PQL/MQL in PROC GLIMMIX in SAS.

```
PROC GLIMMIX METHOD=RSPL;
  CLASS id group occasion;
  MODEL y=group time group*time / DIST=BINOMIAL LINK=LOGIT S;
  RANDOM occasion / SUBJECT=id TYPE=CS RESIDUAL;
```

---

**Table 15.12** Illustrative commands for a marginal logistic regression model, with working correlation specified by introducing randomly varying intercepts, fitted using MQL in PROC GLIMMIX in SAS.

```
PROC GLIMMIX METHOD=RMPL;
  CLASS id group;
  MODEL y=group time group*time / DIST=BINOMIAL LINK=LOGIT S;
  RANDOM INTERCEPT / SUBJECT=id TYPE=UN;
```

---

PROC GLIMMIX <options>;

This statement calls the procedure GLIMMIX in SAS. It includes an option for specifying the method of estimation, using METHOD=<options>. The default is METHOD=RSPL; three other options include METHOD=MSPL, METHOD=RMPL, and METHOD=MMPL (see [Table 15.8](#)). PROC GLIMMIX also includes an implementation of adaptive Gaussian quadrature using the METHOD=QUAD option. There is also an EMPIRICAL=CLASSICAL option that requests that standard errors and test statistics for the fixed effects be based on the classical empirical or “sandwich” estimator; this option is particularly useful when PROC GLIMMIX is used to fit marginal models. PROC GLIMMIX also includes options for various bias-corrected versions of the empirical variance estimator.

CLASS variables;

The CLASS statement is used to define all variables that are to be treated as categorical or factors. By default, this statement will create indicator variables for each factor using a reference group coding, with the last level (where “last” here refers to the level with the largest alphanumeric value) regarded as the reference group. Different sort orders for the CLASS variables can be requested by the ORDER=<option> on the PROC GLIMMIX statement.

MODEL response = <effects> / <options>;

MODEL events/trials = <effects> / <options>;

The MODEL statement specifies the response variable and the covariate effects. The second form of the MODEL statement, with the events/trials syntax, allows the response to be in the form of a ratio of two variables (e.g., counts of the number of successes and the number of trials) and is used for binomial response data. The linear predictor can include both discrete (defined in the CLASS statement) and quantitative (excluded from the CLASS statement) covariates. By default, PROC GLIMMIX includes a column of 1's for the intercept in the model. The SOLUTION or S option is used to produce estimates of the fixed (or covariate) effects.

The option DIST=*keyword* specifies the conditional distribution of the response given the random effects in a GLMM. Alternatively, if PROC GLIMMIX is used to fit a marginal model, the DIST=*keyword* has a somewhat different role. In fitting a marginal model using the GEE approach, the DIST=*keyword* does not specify the distribution for the vector of correlated responses; instead, it specifies the default canonical link function and variance function that happen to be associated with particular exponential family distributions. For example, for a marginal model the option DIST=POISSON does not specify that the response vector (or even its separate components) has a Poisson distribution; instead, it specifies that the mean of the response vector is related to the covariates via a log link function (the canonical link for the Poisson distribution) and the mean and variance of the responses are related by  $\text{Var}(Y) = E(Y) = \mu$  (i.e., the variance function is  $v(\mu) = \mu$ ).

Note that PROC GLIMMIX provides a wide choice of options for the inclusion of a dispersion parameter,  $\phi$ . However, unlike PROC GENMOD, the option for the inclusion of a dispersion parameter does not appear on the MODEL statement; instead, a dispersion parameter can be included by using the “G-side” variant of the RANDOM statement (see below).

The LINK=*keyword* specifies the choice of built-in link function relating the mean response to

the linear predictor. If the LINK=*keyword* is omitted, the default link function is the canonical link function associated with the particular exponential family distribution specified on DIST=*keyword*.

A final option often required when modeling count data is an offset. The OFFSET=*variable* specifies a variable to be used as an offset. For example, in modeling count data the rate is often of more direct interest and the denominator for the counts or “population at risk” (or more specifically, the log of the denominator) can be included as an offset. Note that this variable cannot be a CLASS variable and it should not be included as one of the covariates listed on the MODEL statement.

RANDOM <random-effects> / SUBJECT=subject-effect <options>;

The RANDOM statement is used in two different ways in PROC GLIMMIX: it is used to define the random effects,  $b_i$  (“G-side” effects), and also used to specify the “residual errors” (“R-side” effects).

In a GLMM, the RANDOM statement is used to define the covariates in the design matrix,  $Z_i$ , for the random effects,  $b_i$ . Ordinarily these will be a subset of the covariates included on the MODEL statement. While the MODEL statement is used to define the design matrix for the fixed effects and the RANDOM statement is used to define the design matrix for the random effects, note that an intercept is included by default in the former but not the latter. That is, unlike the MODEL statement, PROC GLIMMIX does not include an intercept in the RANDOM statement by default. However, you can specify INTERCEPT (or INT) as a random effect on the RANDOM statement. The RANDOM statement is also used to specify the structure of the covariance matrix for the random effects, G. The structure of G is specified using the TYPE=option. The random effects can be assumed to be correlated (TYPE=UN) or uncorrelated (TYPE=VC); ordinarily covariance pattern models are not used to account for the covariance among the random effects.

As was noted earlier, unlike PROC MIXED, the GLIMMIX procedure does not have a REPEATED statement for specifying the “residual errors” (“R-side” effects). Instead, it uses the RANDOM statement with either the RESIDUAL option or the reserved \_RESIDUAL\_ keyword in place of the random effects. The use of the RESIDUAL option on the RANDOM statement indicates that the statement refers to the “residual errors” or “R-side” effects. In most respects, the RANDOM statement with the RESIDUAL option is equivalent to the REPEATED statement in PROC MIXED. For example, when using the RANDOM statement with the RESIDUAL option, it is possible to include a variable denoting the “repeated effect.” In the context of longitudinal data, the “repeated effect” is often used to identify the order of the repeated measurements within subjects (see [Table 15.11](#)):

RANDOM <repeated effect> / SUBJECT=subject-effect RESIDUAL;

There is also a \_RESIDUAL\_ keyword that can be used on the RANDOM statement (in place of the <random-effects>):

RANDOM \_RESIDUAL\_ / SUBJECT=subject-effect;

To recap, the statement:

REPEATED / SUBJECT=id TYPE=CS;

in PROC MIXED is equivalent to the following statement in PROC GLIMMIX:

RANDOM \_RESIDUAL\_ / SUBJECT=id TYPE=CS;

Similarly, if it is necessary to include a “repeated effect” that identifies the measurement occasions, then the following statements in PROC MIXED:

CLASS id occasion;

REPEATED occasion / SUBJECT=id TYPE=UN;

are equivalent to the following statements in PROC GLIMMIX:

CLASS id occasion;

RANDOM occasion / SUBJECT=id TYPE=UN RESIDUAL;

Finally, simply adding the `_RESIDUAL_` keyword to the RANDOM statement:

```
RANDOM _RESIDUAL_;
```

specifies a single variance parameter for the “residual errors.” That is, the scale parameter  $\phi$  is no longer regarded as fixed; instead it is estimated from the data (e.g., to allow for overdispersion relative to binomial or Poisson variation).

## 15.7 BASIS OF PQL AND MQL APPROXIMATIONS\*

In Section 15.4 it was noted that the PQL and MQL methods are based on two different approximations to the GLMM. In both cases these approximations yield a standard linear mixed effects model for a “working” or transformed response, denoted by  $Y_{ij}^*$  and  $Y_{ij}^{**}$  for PQL and MQL, respectively. This section provides a brief description of the basis for the approximations used by the PQL and MQL methods. This section is somewhat technical and can be skipped without loss of continuity.

Recall that the model for the conditional mean of a GLMM,

$$(15.6) \quad g\{E(Y_{ij}|b_i)\} = \eta_{ij}^b = X'_{ij}\beta + Z'_{ij}b_i,$$

can also be expressed as

$$(15.7) \quad Y_{ij} = g^{-1}(X'_{ij}\beta + Z'_{ij}b_i) + \epsilon_{ij} = \mu_{ij}^b + \epsilon_{ij},$$

where  $g^{-1}(\cdot)$  is the inverse link function,  $\mu_{ij}^b$  denotes the *conditional* mean of  $Y_{ij}$ , given  $b_i$ , and  $\epsilon_{ij}$  denotes a mean zero random error. Also we assume  $\text{Var}(Y_{ij}|b_i) = \text{Var}(\epsilon_{ij}) = v\{E(Y_{ij}|b_i)\} \phi$ , where  $v(\cdot)$  is a known variance function.

The basis of the PQL method is an approximation to the model given by (15.7) so that the random effects  $b_i$  and the within-subject errors enter into the approximate model in an additive, linear fashion. This type of approximation can be achieved by using what is known as a first-order Taylor series expansion of (15.7) around current estimates  $\hat{\beta}$  and  $\hat{b}_i$ . Taylor series expansions require an understanding of calculus, so we omit many of the technical details here. It will suffice to mention that a first-order Taylor series expansion can be used to express a non-linear function of  $\beta$  and  $b_i$  (e.g.,  $g^{-1}(X'_{ij}\beta + Z'_{ij}b_i)$ ) as a sum or linear combination of  $\beta$  and  $b_i$ . Specifically, this yields the following approximation:

$$\begin{aligned} Y_{ij} &= g^{-1}(X'_{ij}\beta + Z'_{ij}b_i) + \epsilon_{ij} \\ &\approx g^{-1}(X'_{ij}\hat{\beta} + Z'_{ij}\hat{b}_i) + \delta(\hat{\mu}_{ij}^b)\{X'_{ij}(\beta - \hat{\beta}) + Z'_{ij}(b_i - \hat{b}_i)\} + \epsilon_{ij} \\ &= \hat{\mu}_{ij}^b + \delta(\hat{\mu}_{ij}^b)\{X'_{ij}(\beta - \hat{\beta}) + Z'_{ij}(b_i - \hat{b}_i)\} + \epsilon_{ij}, \end{aligned}$$

where  $\delta(\mu_{ij}^b)$  denotes the derivative of  $\mu_{ij}^b$  with respect to the linear predictor,  $\eta_{ij}^b = X'_{ij}\beta + Z'_{ij}b_i$  (in calculus, the “derivative” of a function describes its rate of change, or how quickly that function changes, with respect to some variable). When a canonical link function has been chosen (e.g., logit link function for Bernoulli or log link function for Poisson),  $\delta(\mu_{ij}^b) = v(\mu_{ij}^b)$ , the variance function. The final step is to re-arrange the terms given above so that all of the unknown quantities appear on the right-hand side:

$$\delta^{-1}(\hat{\mu}_{ij}^b)(Y_{ij} - \hat{\mu}_{ij}^b) + X'_{ij}\hat{\beta} + Z'_{ij}\hat{b}_i \approx X'_{ij}\beta + Z'_{ij}b_i + \delta^{-1}(\hat{\mu}_{ij}^b)\epsilon_{ij},$$

where  $\delta^{-1}(\hat{\mu}_{ij}^b) = 1/\delta(\hat{\mu}_{ij}^b)$ . Closer inspection of the right-hand side of the equation above reveals that it conforms to a standard linear mixed effects model, with fixed effects  $\beta$ , random effects  $b_i$ , and within-subject errors  $\delta^{-1}(\hat{\mu}_{ij}^b)\epsilon_{ij}$ . If we denote the left-hand side of the equation by

$$Y_{ij}^* = \delta^{-1}(\hat{\mu}_{ij}^b)(Y_{ij} - \hat{\mu}_{ij}^b) + X'_{ij}\hat{\beta} + Z'_{ij}\hat{b}_i,$$

and let  $\epsilon_{ij}^* = \delta^{-1}(\hat{\mu}_{ij}^b)\epsilon_{ij}$ , the equation can be expressed as

$$Y_{ij}^* \approx X'_{ij}\beta + Z'_{ij}b_i + \epsilon_{ij}^*,$$

where the within-subject errors,  $\epsilon_{ij}^*$ , have mean of zero and variance equal to  $\phi[\delta^{-1}(\hat{\mu}_{ij}^b)]^2 v(\hat{\mu}_{ij}^b)$ . When expressed in this way, the relation to linear mixed effects models is far more transparent. That is, we now have a standard linear mixed effects model for the “working” response (or so-called pseudo-data),  $Y_{ij}^*$ , with fixed effects  $\beta$ , and with the random effects,  $b_i \sim N(0, G)$ , and the within-subject errors,  $\epsilon_{ij}^*$ , entering into the model in an additive, linear fashion. As a result estimation can proceed by iteratively fitting a linear mixed effects model to the updated “working” response  $Y_{ij}^*$  using the 2-

step algorithm outlined in Section 15.2.

Next we illustrate the form of the “working” response  $Y_{ij}^*$  for a mixed effects logistic regression model. Consider the following logistic regression model with randomly varying intercepts (or subject effects):

$$\log \left\{ \frac{\Pr(Y_{ij} = 1|b_i)}{\Pr(Y_{ij} = 0|b_i)} \right\} = \eta_{ij}^b = X'_{ij}\beta + b_i,$$

where the single random effect  $b_i$  is assumed to have a zero mean univariate normal distribution. The model for the conditional mean can be expressed in the following equivalent way:

$$\Pr(Y_{ij} = 1|b_i) = E(Y_{ij}|b_i) = \mu_{ij}^b = \frac{\exp(X'_{ij}\beta + b_i)}{1 + \exp(X'_{ij}\beta + b_i)}.$$

Conditional on  $b_i$ , the  $Y_{ij}$  are assumed to be independent and to have a Bernoulli distribution, with

$$\text{Var}(Y_{ij}|b_i) = v\{E(Y_{ij}|b_i)\} = E(Y_{ij}|b_i)\{1 - E(Y_{ij}|b_i)\},$$

that is,  $v(\mu_{ij}^b) = \mu_{ij}^b(1 - \mu_{ij}^b)$ . Finally, because the canonical link function for the Bernoulli distribution has been adopted,

$$\delta(\mu_{ij}^b) = v(\mu_{ij}^b) = \mu_{ij}^b(1 - \mu_{ij}^b).$$

For current estimates of  $\beta$  and  $b_i$ , the “working” response  $Y_{ij}^*$  is then given by the expression

$$\begin{aligned} Y_{ij}^* &= \delta^{-1}(\hat{\mu}_{ij}^b)(Y_{ij} - \hat{\mu}_{ij}^b) + X'_{ij}\hat{\beta} + \hat{b}_i \\ &= \frac{Y_{ij} - \hat{\mu}_{ij}^b}{\hat{\mu}_{ij}^b(1 - \hat{\mu}_{ij}^b)} + X'_{ij}\hat{\beta} + \hat{b}_i, \end{aligned}$$

where

$$\hat{\mu}_{ij}^b = \frac{\exp(X'_{ij}\hat{\beta} + \hat{b}_i)}{1 + \exp(X'_{ij}\hat{\beta} + \hat{b}_i)}.$$

The steps for deriving the MQL method are very similar except that the approximation is based on a first-order Taylor series expansion of (15.7) around current estimates  $\hat{\beta}$  and  $b_i = 0$  (the mean of the random effects). The expansion yields the following approximation:

$$\begin{aligned} Y_{ij} &= g^{-1}(X'_{ij}\beta + Z'_{ij}b_i) + \epsilon_{ij} \\ &\approx g^{-1}(X'_{ij}\hat{\beta}) + \delta(\hat{\mu}_{ij})(X'_{ij}(\beta - \hat{\beta}) + Z'_{ij}b_i) + \epsilon_{ij} \\ &= \hat{\mu}_{ij} + \delta(\hat{\mu}_{ij})(X'_{ij}(\beta - \hat{\beta}) + Z'_{ij}b_i) + \epsilon_{ij}, \end{aligned}$$

where  $\delta(\mu_{ij})$  denotes the derivative of the *marginal* mean,  $\mu_{ij} = g^{-1}(X'_{ij}\beta)$ , with respect to the *marginal* model linear predictor,  $\eta_{ij} = X'_{ij}\beta$ . When a canonical link function has been chosen (e.g., logit link function for Bernoulli or log link function for Poisson),  $\delta(\mu_{ij}) = v(\hat{\mu}_{ij})$ , the variance function (applied to the marginal mean,  $\mu_{ij}$ ). Note that unlike PQL,  $\mu_{ij}$  depends on the *marginal* linear predictor  $X'_{ij}\beta$  instead of the *conditional* linear predictor  $X'_{ij}\beta + Z'_{ij}b_i$ . As we have already discussed in earlier sections, this has some important implications for interpretation of the resulting MQL estimates of  $\beta$ .

The final step is to re-arrange the terms given above so that all of the unknown quantities appear on the right-hand side:

$$\delta^{-1}(\hat{\mu}_{ij})(Y_{ij} - \hat{\mu}_{ij}) + X'_{ij}\hat{\beta} \approx X'_{ij}\beta + Z'_{ij}b_i + \delta^{-1}(\hat{\mu}_{ij})\epsilon_{ij},$$

where  $\delta^{-1}(\hat{\mu}_{ij}) = 1/\delta(\hat{\mu}_{ij})$ . Closer inspection of the right-hand side of the equation above reveals that it conforms to a standard linear mixed effects model, with fixed effects  $\beta$ , random effects  $b_i$ , and within-subject errors  $\delta^{-1}(\hat{\mu}_{ij})\epsilon_{ij}$ . If we denote the left-hand side of the equation by

$$Y_{ij}^{**} = \delta^{-1}(\hat{\mu}_{ij})(Y_{ij} - \hat{\mu}_{ij}) + X'_{ij}\hat{\beta},$$

and let  $\epsilon_{ij}^{**} = \delta^{-1}(\hat{\mu}_{ij})\epsilon_{ij}$ , the equation can be expressed as

$$Y_{ij}^{**} \approx X'_{ij}\beta + Z'_{ij}b_i + \epsilon_{ij}^{**},$$

where the within-subject errors,  $\epsilon_{ij}^{**}$ , have mean of zero and variance equal to  $\phi[\delta^{-1}(\hat{\mu}_{ij})]^2v(\hat{\mu}_{ij})$ . That is, we now have a standard linear mixed effects model for the “working” response,  $Y_{ij}^{**}$ , with fixed effects  $\beta$ , and with the random effects,  $b_i \sim N(0, G)$ , and the within-subject errors  $\epsilon_{ij}^{**}$  entering into the model in an additive, linear fashion. As with the PQL method, estimation can proceed by iteratively

fitting a linear mixed effects model to the updated “working” response  $Y_{ij}^{**}$ .

Finally, although the PQL and MQL methods both use GLS estimators,

$$\hat{\beta}^* = \left\{ \sum_{i=1}^N (X_i' V_i^{*-1} X_i) \right\}^{-1} \sum_{i=1}^N (X_i' V_i^{*-1} Y_i^*)$$

and

$$\hat{\beta}^{**} = \left\{ \sum_{i=1}^N (X_i' V_i^{**-1} X_i) \right\}^{-1} \sum_{i=1}^N (X_i' V_i^{**-1} Y_i^{**}),$$

respectively, where  $Y_i^*$  and  $Y_i^{**}$  denote the  $n_i \times 1$  vectors of working responses ( $V_i^*$  and  $V_i^{**}$  are the marginal covariance of  $Y_i^*$  and  $Y_i^{**}$  respectively), the former can be shown to be based on a weighted average of *conditional* (i.e., conditional on the random effects,  $b_i$ ) “residuals,”  $(Y_{ij} - \mu_{ij}^b)$ , while the latter is based on a weighted average of *marginal* “residuals,”  $(Y_{ij} - \mu_{ij})$ . Consequently the GLS estimator of  $\beta$  based on  $Y_i^*$  (PQL estimator) yields estimates of the fixed effects in the model for the conditional mean,

$$g\{E(Y_{ij}|b_i)\} = \eta_{ij}^b = X_{ij}' \beta^* + Z_{ij}' b_i.$$

In contrast, the GLS estimator of  $\beta$  based on  $Y_i^{**}$  (MQL estimator) yields estimates of the regression parameters in the following model for the marginal mean,

$$g\{E(Y_{ij})\} = \eta_{ij} = X_{ij}' \beta^{**}.$$

As was discussed in earlier sections of the book,  $\beta^* \neq \beta^{**}$  when a non-linear link function,  $g(\cdot)$ , is adopted (e.g., logit link function). Thus the PQL method should be used when the goal of the analysis is to make subject-specific inferences for the parameters in the GLMM, whereas the MQL method should be used when the goal is to make population-averaged inferences for the regression parameters in a marginal model that assumes the same link function  $g(\cdot)$ .

## **15.8 FURTHER READING**

Breslow (2005) presents a concise, and remarkably clear, review of the statistical literature on approximate methods for estimation and inference for generalized linear mixed effects models.

# Bibliographic Notes

The theoretical foundations for approximate methods for generalized linear mixed effects models can be found in Stiratelli, Laird, and Ware (1984), Schall (1991), Breslow and Clayton (1993), and Wolfinger (1993). For the special case of the logit-normal model, Stiratelli, Laird, and Ware (1984) proposed an approximate method of estimation, based on empirical Bayes ideas, that circumvented the need for numerical integration. Specifically, they avoided the need for numerical integration by approximating the integrand with a simple expansion whose integral has a closed form. The paper by Stiratelli et al. (1984) provided the impetus for the development of a general approach for fitting generalized linear mixed models known as penalized quasi-likelihood (PQL). Schall (1991), Breslow and Clayton (1993), and Wolfinger (1993), motivated PQL as a Laplace approximation to the marginal likelihood for generalized linear mixed models and highlighted the generality of the PQL method.

## Problems

**15.1** In a randomized, double-blind, parallel-group, multicenter study comparing two oral anti-fungal treatments (200 mg/day Itraconazole and 250 mg/day Terbinafine) for toenail infection (De Backer et al., 1998; also see Lesaffre and Spiessens, 2001), patients were evaluated for the degree of onycholysis (the degree of separation of the nail plate from the nail-bed) at baseline (week 0) and at weeks 4, 8, 12, 24, 36, and 48 thereafter. The onycholysis outcome variable is binary (“none or mild” versus “moderate or severe”). The binary outcome was evaluated on 294 patients comprising a total of 1908 measurements. The main objective of the analyses is to compare the effects of the two oral anti-fungal treatments (Itraconazole and Terbinafine) on changes in the probability of the binary onycholysis outcome over the duration of the study.

The raw data are stored in an external file: `toenail.dat`

Each row of the data set contains the following five variables:

ID Y Treatment Month Visit

*Note:* The binary onycholysis outcome variable  $Y$  is coded 0 = none or mild, 1 = moderate or severe. The categorical variable Treatment is coded 1 = Terbinafine, 0 = Itraconazole. The variable Month denotes the exact timing of measurements in months. The variable Visit denotes the visit number (visit numbers 1–7 correspond to scheduled visits at 0, 4, 8, 12, 24, 36, and 48 weeks).

**15.1.1** Consider a generalized linear mixed model, with randomly varying intercepts, for the patient-specific log odds of moderate or severe onycholysis. Using penalized quasi-likelihood, fit a model with linear trends for the log odds over time, with common intercept for the two treatment groups, but different slopes:

$$\text{logit}\{E(Y_{ij}|b_i)\} = (\beta_1 + b_i) + \beta_2 \text{Month}_{ij} + \beta_3 \text{Treatment}_i \times \text{Month}_{ij},$$

where, given  $b_i$ ,  $Y_{ij}$  is assumed to have a Bernoulli distribution. Assume that  $b_i \sim N(0, \sigma_b^2)$ .

**15.1.2** What is the interpretation of the estimate of  $\beta_2$ ?

**15.1.3** What is the interpretation of the estimate of  $\beta_3$ ?

**15.1.4** From the results of the analysis for Problem 15.1.1, what conclusions do you draw about the effect of treatment on changes in the log odds of moderate or severe onycholysis over time? Provide results that support your conclusions.

**15.1.5** Repeat the analysis from Problem 15.1.1, fitting the model using maximum likelihood (ML) with 30 point numerical quadrature.

**15.1.6** Compare and contrast the estimates of  $\beta$  and  $\sigma_b^2$  from fitting the model using ML and PQL. Can you explain why they might differ?

**15.1.7** Repeat the analysis from Problem 15.1.1, fitting the model using marginal quasi-likelihood (MQL) instead of penalized quasi-likelihood (PQL).

**15.1.8** Compare and contrast the estimates of  $\beta$ , especially  $\beta_3$ , from fitting the model using MQL

and PQL. Can you explain why the estimates might differ? Can you provide results from an additional analysis of these data that might support your explanation for the difference between the MQL and PQL estimates?

# *Chapter 16*

## *Contrasting Marginal and Mixed Effects Models*

### **16.1 INTRODUCTION**

In this chapter we compare and contrast marginal and mixed effects models for longitudinal data. There are a number of important distinctions between these two broad classes of models that go beyond simple differences in approaches to accounting for the within-subject association. We emphasize that these two classes of models have different targets of inference and therefore address subtly different questions regarding longitudinal change. In this chapter we highlight the main distinctions and discuss the types of scientific questions addressed by each of the two classes of models.

## 16.2 LINEAR MODELS: A SPECIAL CASE

In Part II we focused on linear models for longitudinal data where the model for the mean response vector can be expressed as

$$(16.1) \quad E(Y_i) = X_i\beta.$$

To account for the positive correlation among the repeated measurements, we described two broad approaches. The first approach is to adopt a covariance pattern model (e.g., autoregressive, Toeplitz) for  $\Sigma_i = \text{Cov}(Y_i)$ . The second approach is to introduce random effects in the model for the mean response,

$$(16.2) \quad E(Y_i|b_i) = X_i\beta + Z_i b_i,$$

where  $b_i$  is a vector of random effects that vary from individual to individual according to a probability distribution (commonly assumed to be multivariate normal). The introduction of random effects induces a random effects covariance structure for  $\Sigma_i$ ,

$$\Sigma_i = Z_i G Z_i' + \sigma^2 I_{n_i},$$

where  $G = \text{Cov}(b_i)$  and  $\sigma^2$  is the variance of the measurement or sampling errors.

When discussing these two different approaches for accounting for the within-subject association, issues concerning the interpretation of  $\beta$  in (16.1) and (16.2) simply did not arise because  $\beta$  has the same interpretation in both models. That is,  $\beta$  describes how the mean response in the study population changes with time and how these changes are related to the covariates. This interpretation of  $\beta$  is transparent when (16.1) is considered. However, the linear mixed effects model given by (16.2) implies the exact same model for the marginal mean response when averaged over the distribution of the random effects. That is,

$$\begin{aligned} E(Y_i) &= E\{E(Y_i|b_i)\} \\ &= E(X_i\beta + Z_i b_i) \\ &= X_i\beta + Z_i E(b_i) \\ &= X_i\beta, \end{aligned}$$

since the vector of random effects has mean zero (i.e.,  $E(b_i) = 0$ ). Thus, in the linear mixed effects model,  $\beta$  can be interpreted *either* as regression coefficients in the model for the conditional mean of  $Y_i$  given  $X_i$  and  $b_i$  (i.e., conditional on  $b_i$ ) or as regression coefficients in the model for the population average mean of  $Y_i$  given  $X_i$ . In the process of taking the expectation or average over the distribution of the random effects, we have implicitly used the property that expectation is a linear operation. This means that the expectation of any linear function of  $b_i$  can be easily evaluated. That is,

$$E(X_i\beta + Z_i b_i) = X_i\beta + Z_i E(b_i),$$

for any constants  $X_i\beta$  and  $Z_i$ . Thus  $\beta$  has the same interpretation in (16.1) and (16.2) because (16.2) is a *linear* mixed effects model (i.e., the right-hand side of (16.2) is a linear function of  $b_i$ ). Put more simply, in the linear mixed effects model  $\beta$  has a marginal interpretation because the average of the linear rates of change over time for individuals is the same as the linear rate of change over time in the population mean response. However, as we will see in the next section, when evaluating expectations of any *non-linear* functions of  $b_i$  we can no longer proceed in this manner. That is, for any non-linear function of  $b_i$ , say  $h(X_i\beta + Z_i b_i)$ ,

$$E\{h(X_i\beta + Z_i b_i)\} \neq h\{X_i\beta + Z_i E(b_i)\}.$$

## 16.3 GENERALIZED LINEAR MODELS

Next we consider the comparison of marginal and mixed effects generalized linear models for longitudinal data. Recall that one of the components in the specification of a generalized linear model is the link function,  $g(\mu_i)$ , which relates the mean of  $Y_i$  to the linear predictor. In the previous discussion of linear models, the link function was the identity function,  $g(\mu_i) = \mu_i$ . For the special case of an identity link function, and hence linear models, the regression parameters  $\beta$  have the same interpretation in both marginal and mixed effects models. In this section we focus on non-linear (or non-identity) link functions and compare the regression parameters in marginal and generalized linear mixed effects models.

Recall that a marginal model for the mean response vector is given by

$$(16.3) \quad g(\mu_i) = g\{E(Y_i)\} = X_i\beta,$$

where  $g(\cdot)$  is an appropriate vector-valued non-linear link function (e.g., logit or log). The regression parameters  $\beta$  in a marginal model have interpretation in terms of changes in the transformed mean response in the study population, and their relation to covariates. For example, when the components of  $Y_i$  are binary and a logit link function is adopted, with

$$\text{logit}(\mu_i) = X_i\beta,$$

the regression parameters have interpretation in terms of changes in the log odds of success in the study population. For any known link function  $g(\cdot)$ , the population means can be expressed in terms of the inverse link function, say  $h(\cdot) = g^{-1}(\cdot)$ ,

$$(16.4) \quad h\{g(\mu_i)\} = \mu_i = E(Y_i) = h(X_i\beta).$$

For example, when the components of  $Y_i$  are binary and a logit link function has been adopted, the model for  $\mu_i$  is

$$\mu_i = h(X_i\beta) = \frac{\exp(X_i\beta)}{1 + \exp(X_i\beta)},$$

where  $h(\cdot)$  is the inverse-logit link function. Whether expressed as (16.3) or (16.4), the regression parameters  $\beta$  in a marginal model describe changes in the transformed population mean response vector,  $\mu_i$ . Note also that  $\mu_i$  depends on the index  $i$  only via fixed and known covariate values.

Next consider generalized linear mixed models where the conditional mean of  $Y_i$ , given a vector of random effects  $b_i$ , is

$$(16.5) \quad g\{E(Y_i|b_i)\} = X_i\beta^* + Z_ib_i,$$

where the random effects  $b_i$  have a distribution with mean zero and covariance matrix  $G$ . Here we denote the fixed effects by  $\beta^*$  to clearly distinguish them from the corresponding regression parameters in the marginal model in (16.3). The regression coefficients  $\beta^*$  have subject-specific interpretations in terms of changes in the transformed mean response for any individual. That is, to interpret any component of  $\beta^*$  we must consider a unit change in the corresponding covariate while holding  $b_i$  fixed. However, the most natural way to hold  $b_i$  fixed at a particular value is to focus on the conditional mean response vector of any given individual. Alternatively, we can compare two individuals who have the same values for  $b_i$  but who differ by a single unit in the corresponding covariate. The former interpretation of any component of  $\beta^*$  is most natural when the covariate is time-varying; the latter interpretation is more natural when the covariate is time-invariant.

Thus, unlike  $\beta$  in marginal models,  $\beta^*$  has interpretation in terms of changes in the transformed mean response for any individual (or the notional comparison of individuals with the same values for  $b_i$ ). The regression coefficients  $\beta^*$  do not describe changes in the transformed mean response in the study population. The implied model for the marginal means can only be obtained by averaging over the distribution of the random effects. This involves taking an expectation of a non-linear function of  $b_i$ ,

$$\begin{aligned}
\mu_i &= E(Y_i) \\
&= E\{E(Y_i|b_i)\} \\
(16.6) \quad &= E\{h(X_i\beta^* + Z_i b_i)\}.
\end{aligned}$$

Note that  $\mu_i$  depends on the index  $i$  only via fixed and known covariate values.

The expression given in (16.6) is the expectation of a non-linear function of  $b_i$ . It must be evaluated from the definition of expectation as a weighted average, weighted according to the distribution of the random effects,

$$(16.7) \quad \mu_i = E(Y_i) = E\{h(X_i\beta^* + Z_i b_i)\} = \int_{-\infty}^{\infty} h(X_i\beta^* + Z_i b_i) f(b_i) db_i,$$

where the integration denotes summation or averaging and  $f(b_i)$  is the probability density function for  $b_i$  (or the “weights” used in the process of averaging). However, the expression for  $E(Y_i)$  given by (16.7) does not, in general, have a closed-form, and moreover, as noted in the previous section,

$$E(Y_i) \neq h(X_i\beta),$$

for any  $\beta$ . For example, consider the logistic regression model with a randomly varying intercept,

$$\text{logit}\{E(Y_i|b_i)\} = X_i\beta^* + b_i,$$

where  $b_i \sim N(0, \sigma_b^2)$ . The implied model for the marginal mean or marginal probability of success is

$$\begin{aligned}
\mu_i &= E(Y_i) \\
&= E\{E(Y_i|b_i)\} \\
&= E\left\{\frac{e^{(X_i\beta^* + b_i)}}{1 + e^{(X_i\beta^* + b_i)}}\right\} \\
&= \int_{-\infty}^{\infty} \frac{e^{(X_i\beta^* + b_i)}}{1 + e^{(X_i\beta^* + b_i)}} \frac{1}{\sqrt{2\pi\sigma_b^2}} e^{-\frac{1}{2}b_i^2/\sigma_b^2} db_i.
\end{aligned}$$

This expression cannot be evaluated in closed-form, and moreover, is not of the logistic regression form,

$$\frac{e^{(X_i\beta)}}{1 + e^{(X_i\beta)}},$$

for any choice of  $\beta$ . That is, marginally (when averaged over the distribution of the random effects) the logit link function is no longer preserved.

For the special case of the logistic regression model with only a single randomly varying intercept (or subject effect),

$$\text{logit}\{E(Y_i|b_i)\} = X_i\beta^* + b_i,$$

where  $b_i \sim N(0, \sigma_b^2)$ , the following approximate relationship holds:

$$\text{logit}\{E(Y_i)\} \approx (1 + k^2\sigma_b^2)^{-\frac{1}{2}} X_i\beta^*,$$

where  $k = \frac{16\sqrt{3}}{15\pi}$ . The derivation of this approximation is not important. What this approximate relationship highlights is how the logistic regression coefficients in the marginal model are attenuated relative to the corresponding fixed effects in the logistic regression model with a randomly varying intercept,

$$\beta \approx \frac{\beta^*}{\sqrt{1 + 0.346\sigma_b^2}},$$

where  $k^2 = 0.346$ . Thus, when  $\text{Var}(b_i) = \sigma_b^2 > 0$ , the marginal logistic regression model parameters,  $\beta$ , are smaller in absolute value than the fixed effects,  $\beta^*$ , in the mixed effects model. In addition the discrepancy between  $\beta$  and  $\beta^*$  increases with increasing  $\sigma_b^2$ . For example, if  $\sigma_b^2 = 3.5$ , then  $\beta^* \approx 1.5 \times \beta$ ; if  $\sigma_b^2 = 9$ , then  $\beta^* \approx 2\beta$ . For the more general logistic regression model with a vector of random effects  $b_i$ , a similar approximate relationship holds and indicates that the marginal logistic regression parameters  $\beta$  are always attenuated toward zero when compared to  $\beta^*$ .

Thus, for non-linear link functions, the fixed effects  $\beta^*$  in generalized linear mixed models are not comparable to the regression parameters  $\beta$  in marginal models. The lack of comparability reflects the distinct targets of inference associated with generalized linear mixed models and marginal models. That is, the fixed effects  $\beta^*$  describe the effects of covariates on changes in an individual's response over time while the regression parameters  $\beta$  describe the effects of covariates on changes in the population mean response over time.

In addition to the special case of an identity function (i.e., linear mixed effects models), there happens to be one exceptional case where  $\beta^*$  and  $\beta$  are almost comparable. When a log link function is adopted, the link function is preserved marginally. In addition the subset of fixed effects  $\beta^*$  for components of  $X_{ij}$  that do not overlap with  $Z_{ij}$  are directly comparable to the corresponding set of marginal regression parameters  $\beta$ . Specifically, for the log-linear regression model with vector of random effects  $b_i$ ,

$$\log\{E(Y_{ij}|b_i)\} = X'_{ij}\beta^* + Z'_{ij}b_i,$$

where  $b_i \sim N(0, G)$ , the implied model for the marginal mean is

$$\begin{aligned}\mu_{ij} &= E\{E(Y_{ij}|b_i)\} \\ &= E\left\{e^{(X'_{ij}\beta^* + Z'_{ij}b_i)}\right\} \\ &= e^{(X'_{ij}\beta^* + \frac{1}{2}Z'_{ij}GZ_{ij})}.\end{aligned}$$

Therefore

$$\log\{E(Y_{ij})\} = X'_{ij}\beta^* + \frac{1}{2}Z'_{ij}GZ_{ij},$$

and  $\beta^*$  differs from  $\beta$  only for those covariates that overlap between  $X_{ij}$  and  $Z_{ij}$ ; for all other covariates, the corresponding components of  $\beta^*$  and  $\beta$  are the same. For example, if the model includes only a single randomly varying intercept (or subject effect),

$$\log\{E(Y_{ij}|b_i)\} = X'_{ij}\beta^* + b_i,$$

then

$$\log\{E(Y_{ij})\} = \log(\mu_{ij}) = X'_{ij}\beta^* + g_{11}/2,$$

where  $g_{11}$  denotes the variance of the randomly varying intercept  $b_i$ . Thus, for the special case of a log link function and a single randomly varying intercept, the fixed effects  $\beta^*$  are directly comparable to the marginal model regression parameters  $\beta$  (with the exception of the intercept).

Finally, although it is possible, in principle, to obtain estimates of the marginal means from a generalized linear mixed effects model, the assumed form for the regression model for the conditional means given  $b_i$  (e.g., logistic or log-linear) no longer holds for the resulting marginal means when averaged over the distribution of the random effects; moreover any misspecification of the model can yield biased estimates of the implied population average means. As a result a set of regression parameters for  $\mu_i$ , describing the dependence of the population mean response on the covariates, is not immediately available from a generalized linear mixed effects model, even after averaging over the distribution of the random effects.<sup>1</sup> The practical consequence is that it is not possible to describe parsimoniously the effects of covariates on the population means in terms of regression coefficients. This may not be so problematic in the setting of a randomized longitudinal clinical trial where the parameter of interest is often a simple difference or contrast of treatment means (or changes in treatment means from baseline). In the latter setting there is only a single covariate of interest (e.g., treatment group) and it is discrete; furthermore a suitable contrast of the *marginal* mean response profiles in the treatment groups can be estimated. However, when one or more of the covariates of interest is quantitative and/or when there are potential confounding variables that need to be controlled for in the analysis, no simple summaries of the effects of covariates on  $\mu_i$  are readily available from generalized linear mixed effects models.

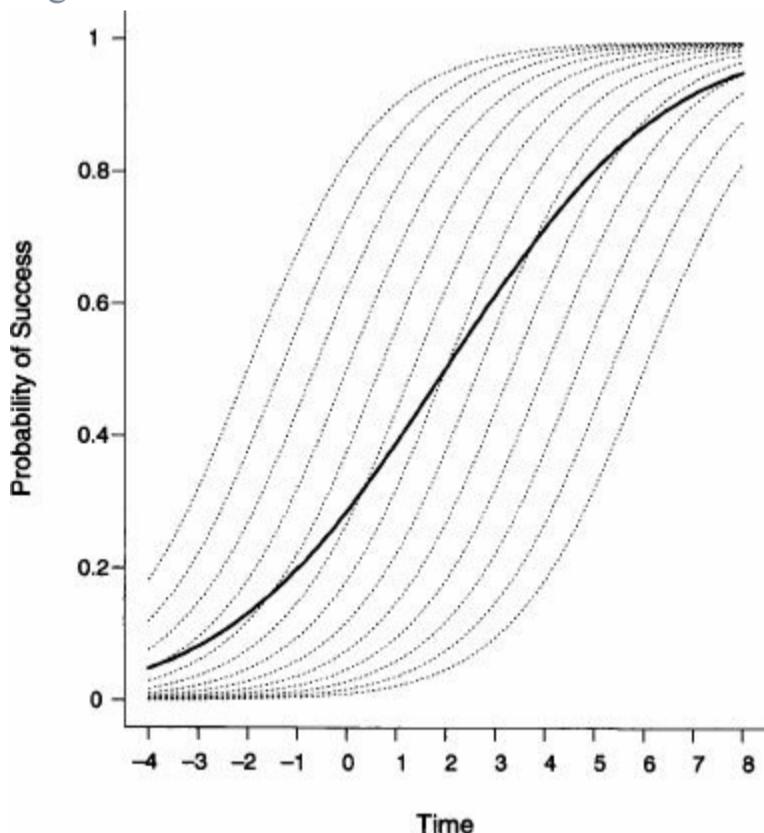
## 16.4 SIMPLE NUMERICAL ILLUSTRATION

To reinforce the distinctions made in the previous section, consider the following simple numerical illustration. Suppose that  $Y_{ij}$  is a vector of binary responses and it is of interest to describe changes in the log odds of success over time. For simplicity we assume that there are no covariates other than the times of measurement. A logistic regression model, with randomly varying intercepts, is given by

$$\text{logit}\{E(Y_{ij}|b_i)\} = \beta_1^* + \beta_2^* t_{ij} + b_i,$$

where  $b_i$  is assumed to have a normal distribution with zero mean and variance  $g_{11} = \text{Var}(b_i)$ . [Figure 16.1](#) displays a plot of  $E(Y_{ij}|b_i)$  versus  $t_{ij}$  for a random sample of  $b_i$  from a normal distribution with zero mean and variance  $g_{11} = 4$  (with  $\beta_1^* = -1.5$  and  $\beta_2^* = 0.75$ ;  $t_{ij} \in [-4, 8]$ ). Also displayed in [Figure 16.1](#) is a plot of the marginal probability of success, averaged over the distribution of  $b_i$ . When the subject-specific logistic curves are compared to the population average curve, it is apparent that the slopes of the former (determined by  $\beta_2^*$ ) are steeper than the slope of the latter. Notice that in the range of probabilities from 0.3 to 0.7, where the logistic curves are approximately linear, the slopes for the subject-specific curves rise faster than the slope for the marginal success probabilities. This reinforces the notion that  $\beta^*$  does not characterize aspects of the population log odds of response, but instead describes changes in the log odds of success for an individual from the population.

**Fig. 16.1** Comparison of conditional probabilities of success (dotted lines) and marginal probability of success (solid line), averaged over the distribution of the random effects.



## 16.5 CASE STUDY

This section highlights aspects of interpretation of the regression coefficients in marginal and generalized linear mixed effects models using safety data from a crossover trial on the disease cerebrovascular deficiency. The variable we analyze is not a trial endpoint per se but rather a potential side effect. In this two-period crossover trial, comparing the effects of active drug to placebo, 67 patients were randomly allocated to the two treatment sequences, with 34 patients receiving placebo → active, and 33 patients receiving active → placebo. The response variable is binary, indicating whether an electrocardiogram (ECG) was abnormal ( $Y = 1$ ) or normal ( $Y = 0$ ). Thus each patient has a bivariate binary response vector,  $Y_i = (Y_{i1}, Y_{i2})'$ , where  $Y_{ij}$  denotes the response for the  $i^{th}$  subject in the  $j^{th}$  period (for  $i = 1, \dots, 67; j = 1, 2$ ). In [Table 16.1](#) the data are summarized in the form of a  $2 \times 4$  contingency table.

**Table 16.1** Data on whether an electrocardiogram (ECG) was normal (0) or abnormal (1) from a two-period crossover trial comparing the effects of active drug to placebo.

Sequence	Response (Period 1, Period 2)			
	(1, 1)	(1, 0)	(0, 1)	(0, 0)
Sequence 1 (P → A)	6	0	6	22
Sequence 2 (A → P)	9	4	2	18

Source: Reprinted with permission from Table 3.1 of Jones and Kenward (1989).

Note: P: Placebo; A: Active drug.

First, we analyze these data using a marginal model. The marginal mean of the response (or probability of an abnormal ECG) is modeled as a logistic function of the covariates, Treatment<sub>ij</sub> (0 = Placebo, 1 = Active drug) and Period<sub>ij</sub> (0 = Period 1, 1 = Period 2),

$$\text{logit}(\mu_{ij}) = \beta_1 + \beta_2 \text{Treatment}_{ij} + \beta_3 \text{Period}_{ij}.$$

The within-subject association between the two responses is modeled in terms of a common log odds ratio,  $\alpha$ ,

$$\text{log OR}(Y_{i1}, Y_{i2}) = \log \left\{ \frac{\Pr(Y_{i1} = 1, Y_{i2} = 1)}{\Pr(Y_{i1} = 1, Y_{i2} = 0)} \frac{\Pr(Y_{i1} = 0, Y_{i2} = 0)}{\Pr(Y_{i1} = 0, Y_{i2} = 1)} \right\} = \alpha.$$

The results, obtained using the GEE approach, are presented in [Table 16.2](#). These results indicate that treatment with the active drug is harmful, increasing the rates of abnormal electrocardiograms. The odds of an abnormal electrocardiogram is 1.77 (or  $e^{0.57}$ ) times higher when treated with active drug versus placebo. The estimate of the within-subject association is  $\hat{\alpha} = 3.56$ , indicating very strong positive association. That is, the odds of an abnormal electrocardiogram at the second occasion is approximately 35 times higher if the electrocardiogram at the first occasion is abnormal rather than normal.

**Table 16.2** Parameter estimates and standard errors from marginal logistic regression model for the ECG data.

Variable	Estimate	SE	Z
Intercept	-1.2433	0.2999	-4.15
Treatment	0.5689	0.2335	2.44
Period	0.2951	0.2319	1.27
log OR( $\alpha$ )	3.5617	0.8148	4.37

Next we analyze these data using a generalized linear mixed model. In particular, we consider the following logistic regression model for the conditional mean of the response, given a random patient effect:

$$\text{logit}\{E(Y_{ij}|b_i)\} = \beta_1^* + \beta_2^* \text{Treatment}_{ij} + \beta_3^* \text{Period}_{ij} + b_i,$$

where the random effect  $b_i$  is assumed to have a normal distribution with zero mean and variance,  $\text{Var}(b_i) = g_{11}$ . In this model each patient is assumed to have some underlying propensity for an abnormal electrocardiogram given by  $b_i$ . Then a patient's odds of an abnormal electrocardiogram is multiplied by a common factor  $e^{\beta_2^*}$  if the patient is treated with the active drug, regardless of that patient's underlying propensity. Thus  $e^{\beta_2^*}$  has interpretation in terms of the ratio of a patient's odds of an abnormal electrocardiogram, when treated with the active drug versus placebo.

The ML estimates of the fixed effects and variance component are presented in [Table 16.3](#). These results also indicate that treatment with the active drug is harmful, increasing the patient-specific rates of abnormal electrocardiograms. In particular, a patient's odds of an abnormal electrocardiogram is 6.4 (or  $e^{1.86}$ ) times higher when treated with active drug than when treated with the placebo. This common treatment effect is the same, regardless of the patient's underlying propensity for an abnormal electrocardiogram. The estimate of the variance of  $b_i$ ,  $\hat{g}_{11} = 24.4$ , indicates that there is very substantial between-patient variability in the propensity for abnormal electrocardiograms, and this is consistent with the very strong within-subject association found in the results from the marginal model reported in [Table 16.2](#). Note, however, that the variance of  $b_i$  is poorly estimated (as evidenced by the very large standard error) since there are only two observations per patient. The estimates of the fixed effects presented in [Table 16.3](#) are all substantially larger than the corresponding estimates in [Table 16.2](#) (approximately three to four times larger), but so too are the standard errors (approximately four times larger). As a result Wald test statistics for null hypotheses concerning the fixed effects are often quite similar in magnitude for the two classes of models. In general, it can be shown that the discrepancy between  $\beta^*$  from the logistic regression model with random subject effect and  $\beta$  from a marginal model (with corresponding fixed effects) increases with the variance of the random subject effect. That is, the greater the underlying heterogeneity among subjects, the greater is the discrepancy between  $\beta^*$  and  $\beta$ , with  $|\beta^*| > |\beta|$  for  $\text{Var}(b_i) > 0$ .

**Table 16.3** Parameter estimates and standard errors from mixed effects logistic regression model for the ECG data.

Variable	Estimate	SE	Z
Intercept	-4.0816	1.6710	-2.44
Treatment	1.8631	0.9269	2.01
Period	1.0376	0.8189	1.27
$g_{11} = \text{Var}(b_i)$	24.4345	18.8474	

Note: ML estimation is based on 100-point adaptive Gaussian quadrature.

Comparison of the two estimates of treatment,  $e^{\hat{\beta}_2} = 1.8$  and  $e^{\hat{\beta}_2} = 6.4$ , from the marginal and mixed effects logistic regression models highlights the distinction between these two analytic approaches. The estimated treatment effect from the marginal model describes how the average rates (expressed in terms of odds) of abnormal electrocardiograms would increase in the study population if patients were treated with the active drug. In contrast, the estimated treatment effect from the mixed effects model describes how the odds of an abnormal electrocardiogram increases for any patient treated with the active drug. As a result the answer to the question "how harmful is the active drug" will depend on whether scientific interest is in its impact on the study population or on an individual drawn at random from that population.

We conclude this section by considering the sensitivity of the results reported in [Tables 16.2](#) and [16.3](#) to sampling zeros and the distributional assumption for the random effects. Note that the first row of [Table 16.1](#) contains a sampling zero. For patients receiving the first sequence of treatments ( $P \rightarrow A$ ), the response pattern (1,0) happens not to have occurred; we refer to this as a "sampling zero" because the response pattern is assumed to be unobserved due to the limited size of the sample (34 patients receiving  $P \rightarrow A$ ). Sampling zeros, an extreme case of sparseness, are known to potentially cause problems for estimation of certain parameters. To assess whether the sampling zero in [Table](#)

[16.1](#) has an inordinate influence on the estimate of the treatment effect, a small constant was added to that cell of the table. Specifically, we added  $\frac{1}{4}$  to the cell with the sampling zero and repeated the two analyses reported in [Tables 16.2](#) and [16.3](#). The estimated treatment effect from the marginal model was  $\hat{\beta}_2 = 0.55$  (SE = 0.236), almost identical to the result found in [Table 16.2](#). The estimated treatment effect from the mixed effects model was  $\hat{\beta}_2^* = 1.64$  (SE = 0.821), somewhat smaller than the estimated treatment effect in [Table 16.3](#); the estimated variance of  $b_i$ ,  $\hat{\sigma}_{11} = 20.7$ , was also smaller. However, from a subject-matter point of view, the substantive conclusions do not change and the difference between the subject-specific and population-averaged effects of treatment is of the same order of magnitude as reported in [Tables 16.2](#) and [16.3](#).

Next, we consider the validity of the distributional assumption for  $b_i$ . In particular, does the large estimated variance of  $b_i$  accurately reflect the true between-patient variability in the risk of an abnormal ECG? Because each patient has two responses, there are only four possible patterns of response: (0, 0), (0, 1), (1, 0) and (1, 1). The “sufficient statistic” for  $b_i$  is each patient’s total number of abnormal ECGs, say  $S_i = \sum_{j=1}^2 Y_{ij}$ ; moreover  $S_i$  can take on only three possible values (0, 1, 2). Therefore for patients randomized to each of the two possible treatment sequences,  $b_i$  can take on only three distinct values ordered from smallest to largest according to  $S_i$ . Furthermore, note from [Table 16.1](#) that the distribution of  $S_i$  is far from symmetric, with 60% of patients having no abnormal ECGs, 18% having 1 abnormal ECG, and 22% having 2 abnormal ECGs. This feature of the distribution of  $S_i$  can only be captured by a normal distribution for the  $b_i$ , centered at zero, that has a relatively large variance. To assess the sensitivity of the results to the normal assumption for  $b_i$ , we estimated the treatment effect using *conditional* maximum likelihood estimation (see Section 14.6). That is, we fitted the following *fixed effects* (or *conditional*) logistic regression model,

$$\text{logit}\{E(Y_{ij}|\alpha_i)\} = \beta_1^* + \beta_2^* \text{Treatment}_{ij} + \beta_3^* \text{Period}_{ij} + \alpha_i,$$

where the  $\alpha_i$  denote *fixed effects* representing time-invariant characteristics of the patients that are not otherwise accounted for by the covariates in the model. The conditional ML estimate of the treatment effect, after we added  $\frac{1}{4}$  to the cell with the sampling zero, was  $\hat{\beta}_2^* = 1.94$  (SE = 1.109), larger than the estimated treatment effect,  $\hat{\beta}_2^* = 1.64$ , from the mixed effects model reported earlier. However, this difference in the estimates of  $\beta_2^*$  is well within the sampling variability of the estimates. A statistical test of the difference between the two estimates for the treatment effect, using the analogue of the “Hausman test” developed for GLMMs by Tchetgen and Coull (2006) (see Section 14.6), yielded  $Z = 0.402$ , ( $p > 0.65$ ). Thus the results presented for these data do not appear to be very sensitive to the normal distribution assumption for  $b_i$ ; however, this cannot be expected in general. In particular, when the number of repeated binary responses is small, and there is a large proportion of subjects with positive (or negative) responses at all occasions, the assumption of a symmetric, normal distribution for the  $b_i$  is questionable.

## 16.6 CONCLUSION

In Chapters 12 through 15 we considered two types of extensions of generalized linear models to longitudinal data: marginal models and generalized linear mixed models. These two quite different analytic approaches arise from different specifications of, or assumptions about, the joint distribution of  $Y_i$  and the source of the correlation among the repeated measures on the same individual. Marginal models merely acknowledge the correlation among repeated measures when estimating the regression parameters; in contrast, generalized linear mixed models provide an explanation for the source of the correlation. Unlike the linear models for continuous responses considered in Part II, with generalized linear models (and non-linear link functions) for discrete responses, different assumptions about the correlation can lead to regression coefficients with quite distinct interpretations.

The basic premise of marginal models is to make inferences about population means, and comparisons of sub-population means, albeit on a transformed scale (e.g., logit or log). The term “marginal” is used to emphasize that the mean response modeled is conditional only on the covariates and not on unobserved random effects or on previous responses. A distinctive feature of marginal models is that the regression models for the mean response and the models for the within-subject association are specified separately. This separation of the model for the mean response from the model for the within-subject association ensures that the marginal model regression coefficients have interpretation that does not depend on the assumptions made about the within-subject association. Specifically, the regression coefficients in marginal models describe the effects of covariates on the population mean response.

In contrast, the basic premise of generalized linear mixed effects models is that there is natural heterogeneity across individuals in the study population in a subset of the regression parameters. That is, a subset of the regression parameters (e.g. intercepts and slopes) is assumed to vary across individuals according to some underlying distribution. But, conditional on the random effects, it is assumed that the repeated measurements for any given individual are independent observations. Generalized linear mixed models extend the conceptual approach of the linear mixed effects model in a very natural way. The correlation among repeated measurements arises from their sharing of common random effects. Unlike the linear mixed effects model, the regression parameters in generalized linear mixed models have subject-specific, but not population-averaged, interpretations. That is, due to the non-linear link functions that are usually adopted for discrete responses, the fixed effects do not describe changes in the mean response in the study population. Instead, they describe how changes in an individual’s mean response are related to within-individual changes in the covariates. As a result generalized linear mixed models are most useful when the scientific objective is to make inferences about individuals rather than the study population. For example, the regression parameters in a logistic mixed effects model describe how the log odds of response changes over time, and how these changes relate to within-individual changes in covariates. Unlike marginal models, they do not compare changes in the log odds of response across sub-populations of individuals defined by values of the covariates. In summary, with generalized linear mixed models, the main focus is on inferences about the individual, while with marginal models, the main focus is on inferences about the study population.

The choice between marginal and generalized linear mixed models for longitudinal data can only be made on subject-matter grounds. We have emphasized the different targets of inference for these two classes of models. For any given longitudinal study, different scientific questions will usually demand different analytic models. For example, a physician considering the potential benefits of a novel treatment for one of her patients might be more interested in the subject-specific effect of treatment. On the other hand, public health researchers or health insurance assessors considering the potential reduction in morbidity or mortality in the population if patients receive the novel treatment would be more interested in the population-averaged effect of treatment. When the answers to both of these questions are of interest, there is no contradiction in reporting estimates of both the subject-specific and population-averaged effects.

In summary, we do not prescribe (or proscribe for that matter) one class of models over another. While there has been much debate in the statistical literature concerning the appropriateness of these two classes of models for analyzing longitudinal data, much of the discussion has generated more heat than light. From a purely probabilistic point of view, generalized linear mixed models might appear to have a distinct advantage over marginal models since the marginal distribution of  $Y_i$ , the target of inference for marginal models, can, in principle, be derived from the generalized linear mixed effects model by averaging over the distribution of the random effects. However, this apparent advantage is somewhat illusory because the induced marginal model does not, in general, retain the same form. For example, consider a logistic regression model with random effects. The implied model for the marginal mean, averaged over the distribution of the random effects, cannot be a logistic regression model; that is, a logistic regression model with random effect is simply not compatible with a logistic regression model for the marginal means (when averaged over the distribution of the random effects). As a result regression coefficients that parsimoniously summarize the covariate effects of interest are not readily available. In addition any misspecification of the generalized linear mixed effects model can yield biased estimates of the implied marginal means. If the goal is to make an inference about the population average mean of  $Y_i$ , a marginal model should be adopted, thereby avoiding the aforementioned problems, the need to correctly specify the conditional distribution of  $Y_i$  given  $b_i$  and the marginal distribution of  $b_i$ , and the computational demands of integrating over the distribution of the random effects. Thus we find ourselves in substantial agreement with Drum and McCullagh (1993, p. 300) when they comment that:

“...the megalomaniacal strategy of fitting a grand unified model, supposedly capable of answering any conceivable question that might be posed, is, in our view, dangerous, unnecessary and counterproductive.”

The answers to different scientific questions concerning longitudinal change will invariably demand that different statistical models have to be applied to the data at hand. In short, one size does not fit all.

## 16.7 FURTHER READING

A useful discussion of the distinct interpretations of the regression parameters in marginal and mixed effects models for binary data can be found in Section 12.2 of Agresti (2002); also see Chapter 7 of Fitzmaurice et al. (2009) for a detailed discussion of the distinct targets of inference for marginal and mixed effects models. Gardiner et al. (2009) compare and contrast the assumptions that underlie fixed effects, random effects, and marginal models.

# Bibliographic Notes

Neuhaus et al. (1991) compare marginal and mixed effects models for analyzing correlated binary data; also, see Zeger et al. (1988), Graubard and Korn (1994), and Section 7.4 of Diggle et al. (2002). Hubbard et al. (2010) give a remarkably lucid description of the assumptions behind, and the type of inference provided by, marginal and mixed effects models in the cluster-correlated data setting.

<sup>1</sup> Note that marginal regression parameters can be estimated directly from models with random effects only when a non-standard relationship is assumed for the conditional mean given the random effects. Alternatively, certain non-standard random effects distributions preserve the same link function relating the marginal and conditional means to the covariates. This class of models is known as *marginalized* mixed effects models.

## *Part IV*

### *Missing Data and Dropout*

# *Chapter 17*

## *Missing Data and Dropout: Overview of Concepts and Methods*

### **17.1 INTRODUCTION**

Missing data are a common and challenging problem in the analysis of longitudinal data. Although most longitudinal studies are designed to collect data on every individual in the sample at each time of follow-up, many studies have some missing observations. In longitudinal studies in the health sciences, missing data are the rule, not the exception. With longitudinal studies problems of missing data are far more acute than in cross-sectional studies, since non-response can occur at any occasion. An individual's response can be missing at one follow-up time and then be measured at a later follow-up time, resulting in a large number of distinct missingness patterns. Alternatively, longitudinal studies often suffer from the problem of attrition or "dropout"; that is, some individuals "drop out" or withdraw from the study before its intended completion. In either case the term "missing data" is used to indicate that an intended measurement could not be obtained.

Missing data have three important implications for longitudinal analysis. First, when data are missing, the data set is necessarily unbalanced over time since not all individuals have the same number of repeated measurements at a common set of occasions. As a result missing data create complications for methods of analysis that require balanced data. However, imbalance created by missingness is not a concern for the regression methods described in Parts II and III. Second, when there are missing data, there must necessarily be some loss of information. That is, missing data reduce the precision with which changes in the mean response over time can be estimated. Not surprisingly, the reduction in precision is directly related to the amount of missing data; that is, the greater the amount of missing data, the greater is the decrease in precision. The loss of precision can also depend to some extent on the method of analysis; for example, analyses restricted to subjects with complete data will generally be less efficient than methods that use all available data. The location of the missing data (e.g., missingness spread sporadically over many subjects, or concentrated at a specific set of time points in a few subjects), and how highly correlated the missing data are with the observed data, will also affect loss of precision. Finally, under certain circumstances missing data can introduce bias and thereby lead to misleading inferences about changes in the mean response. It is this last factor, the potential for serious bias, that complicates the analysis of partially missing longitudinal data. As a result the reasons for any missing data, often referred to as the *missing data mechanism*, must be carefully considered.

This is the basis for an important theme that will be emphasized throughout this chapter: when data are missing, we must carefully consider why they are missing. Some types of missing data are relatively benign and do not complicate the analysis; others are not and can potentially introduce bias in the estimates of the regression parameters. The following two examples of partially missing longitudinal data will help illustrate this point.

The first example is from the Six Cities Study of Air Pollution and Health, discussed in Sections 8.8 and 9.6. This was a longitudinal study designed to characterize lung function growth as measured by changes in pulmonary function in children and adolescents. Most of the children were enrolled in the first or second grade (between the ages of six and seven) and measurements were obtained annually until graduation from high school or loss to follow-up. At each annual examination, spirometry, the measurement of pulmonary function, was performed. Due to late entry into the study and loss to follow-up or attrition, the number of measurements of pulmonary function of study

children varied from a minimum of 1 to a maximum of 12. The major reason for late entry or attrition was moving in or out of the school district. Let us focus on this main reason for missing data. If a child changed school district because of employment relocation by her parents, then the missing data mechanism can be thought of as unrelated to the child's pulmonary function. On the other hand, if a child moved out of the school district because she developed respiratory problems (e.g., relocating to an area with either better air quality or improved access to health care), then missingness is related to the child's pulmonary function.

The second example is from the Muscatine Coronary Risk Factor (MCRF) study, introduced in Section 1.3 and analyzed in Section 13.4. This was a longitudinal survey of school-age children in Muscatine, Iowa, examining the development and persistence of risk factors for coronary disease. In the MCRF study, weight and height measurements of five cohorts of children, initially aged 5–7, 7–9, 9–11, 11–13, and 13–15 years, were obtained biennially from 1977 to 1981. On the basis of a comparison of their weight to age–gender-specific norms, children were classified as obese or not obese. The study protocol required parental consent prior to each measurement. One objective of the MCRF study was to determine whether the risk of obesity increased with age and whether patterns of change in obesity were the same for boys and girls. Although each child was eligible to participate in all three surveys, there was a substantial amount of missing data on obesity, with less than 40% of the children providing complete data at all three measurement occasions. The two main reasons for missing data were: (1) failure to obtain consent and (2) the child's absence from school on the day of examination. Let us focus on these two reasons for missing data. Suppose that parents of children who were obese were either more or less likely to sign the consent form than parents of children who were not obese. For example, parents of children who were obese might have been more likely to sign the consent form due to concerns about the adverse health effects of obesity; conversely, they might have been less likely to sign the consent form due to concerns that participation in the study could be a source of embarrassment for their children. In either case the reason for missing data on weight and height is related to the obesity status of the child. Similarly missingness is related to the obesity status of the child if children who were obese were more likely to be absent on the day of examination (e.g., due to embarrassment about being overweight). On the other hand, suppose that a child was absent on the day of examination because of employment relocation by her parents (completely unrelated to the health of the child). Then missingness does not depend on the child's obesity status.

These two examples show that there can be more than a single cause of missing data and that reasons for missing data may or may not be related to the outcome of interest. When data are missing for reasons unrelated to the outcome of interest, the impact of missing data is relatively benign and does not complicate the analysis. When it is related to the outcome, somewhat greater care is required because there is potential for bias when individuals with missing data differ in important ways from those with complete data.

In this chapter we review three general models for missing data that differ in terms of assumptions concerning whether missingness is related to observed and unobserved responses. We also discuss the implicit assumptions about missing data that underlie the methods for longitudinal analysis described in Parts II and III. We illustrate the main distinctions between the three general models for missing data for the common problem of dropout. Finally, we provide an overview of some alternative methods for handling dropout in longitudinal studies. A more in-depth discussion and application of two important methods for handling missing data, multiple imputation and inverse probability weighted methods, can be found in Chapter 18.

## 17.2 HIERARCHY OF MISSING DATA MECHANISMS

To obtain valid inferences from partially missing longitudinal data, we must consider the nature of the “missing data mechanism.” Ordinarily the missing data mechanism is not under the control of the investigators and often is not well understood. Instead, assumptions are made about the missing data mechanism and the validity of the analysis depends on whether these assumptions hold. When reporting the results of a longitudinal analysis, it is important to be explicit about the assumptions made regarding the reasons for missing data.

The missing data mechanism can be thought of as a probability model for the distribution of a set of response indicator variables. These response indicator variables take the value 1 when an intended measurement of the response is obtained and the value 0 otherwise. For example, suppose that the design of the study calls for  $n$  measurements per subject. That is, we intend to take  $n$  repeated measures of the response variable on the same individual. A subject with a *complete* set of responses has an  $n \times 1$  response vector denoted by

$$Y_i = (Y_{i1}, Y_{i2}, \dots, Y_{in})'$$

this is a slight abuse of the notation adopted in previous chapters where  $Y_i$  contained not the possible set of observations but the values actually observed. Because of missing data, some of the components of  $Y_i$  are not observed for at least some individuals. We let  $R_i$  be an  $n \times 1$  vector of response indicators

$$R_i = (R_{i1}, R_{i2}, \dots, R_{in})'$$

with  $R_{ij} = 1$  if  $Y_{ij}$  is observed and  $R_{ij} = 0$  if  $Y_{ij}$  is missing. In addition associated with  $Y_i$  is an  $n \times p$  matrix of covariates,  $X_i$ . We do not consider missingness in the covariates; that is, we assume that any time-varying covariates are fixed by the study design. Given  $R_i$ , the complete set of responses,  $Y_i = (Y_{i1}, \dots, Y_{in})'$ , can be partitioned into two components  $Y_i^O$  and  $Y_i^M$ , corresponding to those responses that are observed and missing, respectively. That is,  $Y_i^O$  denotes the vector of *observed* responses on the  $i^{th}$  subject, and  $Y_i^M$  denotes the complementary set of responses that are missing. The random vector  $R_i$  is recorded for all individuals. Also, given  $R_i$ , the target population of interest can be divided or stratified into a number of distinct sub-populations defined by the missing data patterns (including the sub-population of “completers”). Thus  $R_i$  can also be thought of as a stratification variable that divides the target population into a number of sub-populations. This is illustrated in [Table 17.1](#), where the first response,  $Y_1$ , perhaps denoting a baseline response, is fully observed, but  $Y_2, \dots, Y_n$  are missing intermittently.

**Table 17.1** Schematic representation of  $R$ , the vector of response indicators, as a stratification variable.

Response Indicators						Response Vector <sup>a</sup>					
$R_1$	$R_2$	$R_3$	$R_4$	...	$R_n$	$Y_1$	$Y_2$	$Y_3$	$Y_4$	...	$Y_n$
1	1	1	1	...	1	$y_1$	$y_2$	$y_3$	$y_4$	...	$y_n$
1	0	1	1	...	1	$y_1$	*	$y_3$	$y_4$	...	$y_n$
1	1	0	1	...	1	$y_1$	$y_2$	*	$y_4$	...	$y_n$
1	1	1	0	...	1	$y_1$	$y_2$	$y_3$	*	...	$y_n$
:	:	:	:	...	:	:	:	:	:	...	:
1	0	0	0	...	1	$y_1$	*	*	*	...	$y_n$
1	0	0	0	...	0	$y_1$	*	*	*	...	*

<sup>a</sup>The \* denotes missing value.

A hierarchy of three different types of missing data mechanisms can be distinguished by considering how  $R_i$  is related to  $Y_i$ :

**1. Missing Completely at Random (MCAR),**

**2. Missing at Random (MAR), and**

**3. Not Missing at Random (NMAR).**

The hierarchy of missing data mechanisms is useful because the type of missing data mechanism determines the appropriateness of different methods of analyses, for example, maximum likelihood, generalized least squares (GLS), and GEE. We discuss this topic later in the chapter. However, the nomenclature is not intuitive and leads to much confusion among statisticians and practitioners alike. A major objective of this chapter is to explain these mechanisms in a more intuitive manner so that the reader gains a better appreciation for their usage.

Much of the remainder of this section is devoted to a detailed explanation, with concrete examples, of this classification of missing data mechanisms in the context of longitudinal studies. We begin each description with the formal definition expressed as conditions on the probability distribution of the response indicators,  $R_i$ . We then provide some examples and explain the consequences of each type of missingness for the distribution of the observed data. Once the main distinctions are understood, we can describe the implicit assumptions about the missing data mechanism made by different methods for analyzing longitudinal data.

# Missing Completely at Random (MCAR)

Data are said to be missing completely at random (MCAR) when the probability that responses are missing is unrelated to either the specific values that, in principle, should have been obtained or the set of observed responses. That is, longitudinal data are MCAR when  $R_i$  is independent of both  $Y_i^O$  and  $Y_i^M$ , the observed and unobserved components of  $Y_i$ , respectively. To better understand this missing data mechanism, consider the bivariate case where  $Y_i = (Y_{i1}, Y_{i2})'$ ,  $Y_{i1}$  is assumed to be fully observed and  $Y_{i2}$  is sometimes missing. In that case we require only a single response indicator, with  $R_{i2} = 1$  if  $Y_{i2}$  is observed and  $R_{i2} = 0$  if  $Y_{i2}$  is missing. If  $Y_{i2}$  is MCAR, then

$$\Pr(R_{i2} = 1|Y_{i1}, Y_{i2}, X_i) = \Pr(R_{i2} = 1|X_i),$$

and the probability that  $Y_{i2}$  is missing does not depend on the observed value of  $Y_{i1}$  or the value of  $Y_{i2}$  that, in principle, should have been obtained. Missingness in  $Y_{i2}$  is simply the result of a chance mechanism that does not depend on observed or unobserved components of  $Y_i$ .

An example where partially missing longitudinal data are MCAR is the “rotating panel” study design. In this study design, commonly used in health surveys to reduce response burden, individuals rotate in and out of the study after providing a pre-determined number of repeated measures. However, the number and timing of the measurements are determined by design. The decision about whether to obtain a measurement on an individual at any specific occasion is made a priori by the investigators and is not related to the vector of responses; that is,  $R_i$  is unrelated to  $Y_i$ . In this example the missing data mechanism is under the control of the investigators and is well understood. Another example where missing data are MCAR is in the Six Cities Study of Air Pollution and Health, when children changed school district because of employment relocation by their parents (for reasons completely unrelated to the health of their children). Here the reason for missing data is unrelated to the children’s pulmonary function.

In the definition of MCAR given above, missingness can depend on the covariates,  $X_i$ . This raises a subtle, but important, point. Under MCAR, the response vector  $Y_i$  is conditionally independent of  $R_i$ , given the covariates  $X_i$ . However, this conditional independence of  $Y_i$  and  $R_i$  may not hold when conditioning on only a subset of the covariates. This has the following important implication. When an analysis is based on a subset of  $X_i$  that excludes a covariate that is predictive of  $R_i$ , the missing data can no longer be considered MCAR. For example, in a clinical trial missingness may be related to side effects of the treatments. However, side effects is a covariate that would not ordinarily be included in the analysis model that evaluates treatment effects. If side effects is excluded from the analysis model, the missing data can no longer be considered MCAR; in Chapter 18 we discuss how this more complex case can be handled. When

$$(17.1) \quad \Pr(R_i|Y_i^O, Y_i^M, X_i) = \Pr(R_i|X_i),$$

the data are said to have *covariate-dependent* missingness and use of the term MCAR is sometimes restricted to the case where

$$(17.2) \quad \Pr(R_i|Y_i^O, Y_i^M, X_i) = \Pr(R_i).$$

In our discussion of missing data mechanisms we do not make this subtle distinction. Instead, we define MCAR using (17.1) and simply assume that  $X_i$  in (17.1) contains all relevant covariates for predicting both  $Y_i$  and  $R_i$ .

The essential feature of MCAR is that the observed data can be thought of as a random sample of the complete data. The result is that the moments (e.g., means, variances, and covariances) and, indeed, the distribution of the observed data do not differ from the corresponding moments or distribution of the complete data. Thus, “completers” can be regarded as a random sample from the target population, albeit with a smaller sample size than intended. This has important implications for the analysis of longitudinal data restricted to subjects with complete response vectors. In particular, any method of analysis that yields valid inferences in the absence of missing data will also yield valid inferences when missing data are MCAR and the analysis is restricted to the

“completers.” The latter is often referred to as a “complete-case” analysis. With data MCAR it is legitimate (although possibly wasteful) to remove subjects with any missing data from the analysis since we can regard them as randomly chosen without regard to their data values. This feature of MCAR allows one to do a complete-case analysis without being concerned that the results might be biased by excluding those with missing data.

A similar result holds for subjects with some missing data. The responses actually obtained,  $Y_i^O$ , have the same distribution as the corresponding elements of the completers. As a result all available data can be used to give valid estimates of moments such as means, variances, and covariances. For example, if we modify our bivariate example to allow data to be MCAR either at time 1 or time 2, then subjects with only one observation can be used along with the complete cases to estimate means and variances; only the complete cases can here be used to estimate the covariance. In longitudinal designs with more observations per subject, the observed cases with at least two observations contribute to covariance estimation. As a result methods for longitudinal analysis that incorporate all of the available observations will yield valid inferences when missing data are MCAR. This includes all of the methods that were discussed in Parts II and III of this book.

These properties of MCAR follow directly from the definitions in [\(17.1\)](#) and [\(17.2\)](#). They can be used to show that when the missing data mechanism is MCAR, the distribution of  $Y_i$  (given  $X_i$ ) is the same in each of the distinct sub-populations defined by the missing data patterns (including the sub-population of “completers” or subjects with no missing responses). It also implies that these distributions coincide with the distribution of  $Y_i$  (given  $X_i$ ) in the target population of interest. Moreover, when the missing data mechanism is MCAR, the distribution of the observed components  $Y_i^O$  for any distinct sub-population defined by the missing data patterns coincides with the distribution of the same components of  $Y_i$  in the target population.

Finally, we note that with MCAR, the distribution of  $Y_i^M$  for any distinct sub-population defined by the missing data patterns coincides with the distribution of the same components of  $Y_i$  for the “completers.” For example, in the bivariate case, MCAR implies that the distribution of  $Y_{i2}$  for those missing  $Y_{i1}$  is the same as the distribution of  $Y_{i2}$  for those with no missing responses.

# Missing at Random (MAR)

Data are said to be missing at random (MAR) when the probability that responses are missing depends on the set of observed responses but is conditionally unrelated to the specific missing values that, in principle, should have been obtained. Specifically, longitudinal data are MAR when  $R_i$  is conditionally independent of  $Y_i^M$ , given  $Y_i^O$ ,

$$(17.3) \Pr(R_i|Y_i^O, Y_i^M, X_i) = \Pr(R_i|Y_i^O, X_i).$$

Let us return to the bivariate case where  $Y_i = (Y_{i1}, Y_{i2})'$ ,  $Y_{i1}$  is fully observed, and  $Y_{i2}$  is sometimes missing. If  $Y_{i2}$  is MAR, then

$$\Pr(R_{i2} = 1|Y_{i1}, Y_{i2}, X_i) = \Pr(R_{i2} = 1|Y_{i1}, X_i),$$

and the probability that  $Y_{i2}$  is missing depends on the observed value of  $Y_{i1}$ . However, given  $Y_{i1}$ , the probability that  $Y_{i2}$  is missing does not depend on the value of  $Y_{i2}$  that should have been obtained. Put another way, if subjects are stratified on the basis of similar values for  $Y_{i1}$ , missingness in  $Y_{i2}$  within strata is simply the result of a chance mechanism that does not depend on values of  $Y_{i2}$ .

Longitudinal data that are MAR might arise when a study protocol requires that a subject be removed from the study whenever the value of an outcome variable falls outside of a certain clinical range of values. In that case missingness in  $Y_i$  is under the control of the investigator and is related to observed components of  $Y_i$  only.

Another example where missing data are MAR is in the Six Cities Study of Air Pollution and Health, when children moved out of the school district because they developed respiratory problems. If the decision to relocate could be predicted based only on the recorded history of pulmonary function measurements (i.e., the observed components of  $Y_i$  only), then the missing data are MAR. However, the MAR assumption would not hold if the decision to relocate was based on some extraneous variable, unavailable to the investigators, that was predictive of the future but unobserved, pulmonary function measurements.

Because the missing data mechanism now depends on  $Y_i^O$ , the distribution of  $Y_i$  in each of the distinct sub-populations defined by the missing data patterns is not the same as the distribution of  $Y_i$  in the target population. This has important consequences for analysis. One is that a “complete-case” analysis is not valid and can produce biased estimates of change in the mean response over time. Furthermore the distribution of  $Y_i^O$ , the observed components of  $Y_i$ , in these sub-populations does not coincide with the distribution of the same components of  $Y_i$  in the target population. Therefore the sample means, variances, and covariances based on the available data are biased estimates of the corresponding parameters in the target population. This feature of MAR will be illustrated in the context of dropout in Section 17.4.

With MAR, the observed data cannot be viewed as a random sample of the complete data, but there is an important implication for the distribution of the missing data. The distribution of each individual's missing values,  $Y_i^M$ , conditioned on the observed values,  $Y_i^O$ , is the same as the conditional distribution of the corresponding observations among the complete cases, conditional on those complete cases having the same values as  $Y_i^O$ . In other words, if we stratify on values of  $Y_i^O$ , the distribution of  $Y_i^M$  is the same as the distribution of the corresponding observations in the complete-case and target populations. As a result missing values can be validly predicted using the observed data and a model for the joint distribution. However, the validity of the predictions of the missing values rests on having correctly specified both the model for the mean and the model for the covariance (when the responses have a multivariate normal distribution). The model for the covariance must be correctly specified because conditional moments (e.g., conditional means) depend on both the mean response vector and the covariance.

For example, if the responses are assumed to have a multivariate normal distribution, then predictions of the missing values are based on the conditional mean of  $Y_i^M$ , given  $Y_i^O$ . Using well-

known properties of the multivariate normal distribution, the conditional mean of  $Y_i^M$ , given  $Y_i^O$ , can be expressed as

$$E(Y_i^M|Y_i^O) = \mu_i^M + \Sigma_i^{MO} \Sigma_i^{O^{-1}} (Y_i^O - \mu_i^O),$$

where  $\mu_i^M$  and  $\mu_i^O$  denote those components of the mean response vector corresponding to  $Y_i^M$  and  $Y_i^O$ , and  $\Sigma_i^O$  and  $\Sigma_i^{MO}$  denote those components of the covariance matrix corresponding to the covariance among the elements of  $Y_i^O$  and the covariance between  $Y_i^M$  and  $Y_i^O$ . The important aspect of the expression given above is the dependence of the prediction of  $Y_i^M$  on both the mean response vector

$$\mu_i = \begin{pmatrix} \mu_i^O \\ \mu_i^M \end{pmatrix},$$

and the covariance matrix

$$\Sigma_i = \begin{pmatrix} \Sigma_i^O & \Sigma_i^{OM} \\ \Sigma_i^{MO} & \Sigma_i^M \end{pmatrix}.$$

When missing data are MAR, we must correctly model the entire joint distribution of  $Y_i, f(Y_i|X_i)$  (e.g., both the mean and covariance when  $Y_i$  is assumed to have a multivariate normal distribution) to obtain valid estimates of  $\beta$  (and  $\Sigma_i$ ).

With MAR, the missing values can be predicted using the observed data and a model for the joint distribution of  $Y_i$ . But one does not need to use the model for  $\Pr(R_i|Y_i^O, X_i)$  as a function of  $X_i$  and  $Y_i^O$ , only a model for  $Y_i$  given  $X_i$ . Since MCAR is a special case of MAR, the same is also true of MCAR; namely one does not need to use the model for  $\Pr(R_i|Y_i, X_i)$  to obtain valid likelihood-based inferences, only a model for  $f(Y_i|X_i)$ . Notice that not using  $\Pr(R_i|Y_i, X_i)$  in the analysis has the important implication that we do not need to even posit a specific model for  $\Pr(R_i|Y_i, X_i)$  other than to say it does not depend on the missing observations. Since it is common to use a model for  $f(Y_i|X_i)$ , valid likelihood-based analyses can be obtained with MAR or MCAR data with no extra assumptions, other than the general statement of MCAR or MAR. For this reason MCAR and MAR are often referred to as *ignorable* mechanisms. A caveat concerning this use of the term *ignorable*: when data are MAR, it emphatically does not mean we can ignore the missing data problem and use any complete-case or available-data analysis we desire. Instead, the ignorability refers to the fact that once we establish that  $\Pr(R_i|Y_i, X_i)$  does not depend on missing observations, we can ignore  $\Pr(R_i|Y_i, X_i)$  and obtain a valid likelihood-based analysis provided that we have a correct model for  $f(Y_i|X_i)$ . That is, ML estimation of  $\beta$  in the linear models discussed in Part II is valid when data are MAR provided that the multivariate normal distribution has been correctly specified.

In contrast to a full-likelihood analysis, standard applications of generalized least squares (GLS) that only require a model for the mean response, but do not assume a multivariate normal distribution for the response vector, no longer provide valid estimates of  $\beta$ . That is, when data are MAR, GLS based only on the means, variances, and covariances of the available data can yield biased estimates of  $\beta$ . This is because the sample means, variances, and covariances based on the available data (or based on the “completers”) are biased estimates of the corresponding parameters in the target population. Moreover, with GLS, the means, variances, and covariances may possibly be misspecified. In a similar way the generalized linear mixed effects models (GLMMs) described in Part III fully specify the joint distributions of both the vector of responses and the vector of random effects. As a result conventional likelihood-based analyses of the incomplete data using GLMMs yield valid inferences when data are missing at random (MAR), provided that the likelihood has been correctly specified; this property, however, does not extend to approximate methods such as PQL. This is in contrast to the GEE methods for analyzing discrete longitudinal data described in Part III; the GEE methods require a model for the mean response but do not specify the multivariate joint distribution for the response vector. As a result standard GEE methods do not provide valid estimates of the regression parameters when data are MAR but not MCAR. However, both the GLS

and GEE estimators of  $\beta$  can be adapted to provide a valid analysis by explicitly modeling  $\Pr(R_i | Y_i, X_i)$  and weighting the analysis accordingly; the intuition for “weighted methods” is discussed in Section 17.5 and weighting methods are described in greater detail in Chapter 18.

The subtle distinction between MCAR and MAR is often not well understood. We find that statisticians and empirical researchers regularly confuse the definition of MAR with MCAR (and admittedly, the choice of terminology has not helped matters). As we will see in the next section, the distinction between MAR and MCAR has very important implications for the validity of different methods of analysis of longitudinal data. The MAR assumption is far less restrictive on  $\Pr(R_i)$  than MCAR and may be considered to be a more plausible assumption about missing data in many applications. Of note, although the MAR assumption is less restrictive in the sense of restrictions on  $\Pr(R_i)$ , it can be considered more restrictive in terms of what methods of analyses are appropriate. In our view, the MAR assumption should be the default assumption for the analysis of partially missing longitudinal data unless there is a strong and compelling rationale to support the MCAR assumption.

# Not Missing at Random

The third type of missing data mechanism is referred to as *not missing at random* (NMAR). Missing data are said to be NMAR when the probability that responses are missing is related to the specific values that should have been obtained. That is, the conditional distribution of  $R_i$ , given  $Y_i^O$ , is related to  $Y_i^M$ , and

$$\Pr(R_i|Y_i^O, Y_i^M, X_i)$$

depends on at least some elements of  $Y_i^M$ . Let us return to the bivariate case where  $Y_i = (Y_{i1}, Y_{i2})'$ ,  $Y_{i1}$  is fully observed, and  $Y_{i2}$  is sometimes missing. If missingness in  $Y_{i2}$  is NMAR, then

$$\Pr(R_{i2} = 1|Y_{i1}, Y_{i2}, X_i)$$

depends on the potentially unobserved value of  $Y_{i2}^2$ . An NMAR mechanism is often referred to as *nonignorable* missingness. The term *nonignorable* refers to the fact that the missing data mechanism cannot be ignored when the goal is to make inferences about the distribution of the complete longitudinal responses.

An example where longitudinal data are NMAR arises when the outcome variable is a measure of “quality-of-life” and subjects fail to complete the instrument or questionnaire on occasions when their quality-of-life is compromised. Another example where missing data are NMAR is in the Muscatine Coronary Risk Factor (MCRF) study, when parents of children who were obese were either more or less likely to sign the consent form than parents of children who were not obese. In that case missingness on weight and height is related to the obesity status of the child.

Sometimes the term *informative* is used to describe data that are NMAR; missingness is informative in the sense that the missingness (i.e., a component of  $R_i$  is equal to 0) informs us about the distribution of the missing observations. Specifically, the distribution of  $Y_i^M$ , conditional on  $Y_i^O$ , is not the same as that in the “completers” or in the target population, but rather, the distribution of  $Y_i^M$  depends on  $Y_i^O$  and on  $\Pr(R_i|Y_i, X_i)$ . Thus the model assumed for  $\Pr(R_i|Y_i, X_i)$  is crucial; it must be included in the analysis, and the specific model chosen can drive the results of the analysis.

## 17.3 IMPLICATIONS FOR LONGITUDINAL ANALYSIS

The statistical methods for analyzing longitudinal data presented in earlier chapters of the book can accommodate incomplete data. However, valid inferences from partially missing longitudinal data require assumptions about the missing data mechanism. In this section we summarize the key assumptions about missing data required for valid inferences when applying the techniques described in Parts II and III.

When the missing data mechanism is MCAR, individuals with missing data are a random subset of the sample. In this case the observed values of the responses are a random subsample of all values of the responses, and no bias will arise with almost any method of analysis of the data (either the available data or the data on the “completers” only). In particular, all of the methods discussed in Parts II and III will yield valid estimates of mean response trends (and within-subject associations) if the missing data can be assumed to be MCAR.

When the missing data mechanism is MAR, individuals with missing data are no longer a random subset of the sample. Only when stratified on their observed outcomes (i.e., conditional on  $Y_i^O$ ) can they be considered a random subset of the sample belonging to that stratum. As a result the observed values are not necessarily a random subsample of the responses. In particular, the distribution of  $Y_i^O$ , the observed components of  $Y_i$ , differs from the distribution of the same components of  $Y_i$  in the target population. This implies that based on the available observations the sample means at each occasion (and the covariances) provide biased estimates of the means (and covariances) in the target population. Similarly analyses restricted to the data from the completers also yield biased estimates of the means (and covariances). When missing data are MAR, but not MCAR, complete-case methods and standard GEE methods based on all of the available observations yield biased estimates of mean response trends. In contrast, likelihood-based methods that correctly specify the entire joint distribution of the responses yield valid estimates when missing data are MAR. However, there is a subtle, but important, proviso: the models for both the mean response and the within-subject association must be correctly specified. Thus, when missing data are MAR, the likelihood-based methods discussed in Part II provide valid inferences about changes in the mean response over time provided that the covariance matrix has been correctly modeled. Similarly the methods discussed in Chapter 14 provide valid estimates of the fixed effects provided that the random effects structure has been correctly specified. In summary, when missing data are MAR, but not MCAR, inferences about the mean response are sensitive to any form of misspecification of the joint distribution of the vector of responses. Accordingly, if longitudinal data are incomplete, somewhat greater care must be exercised when modeling the within-subject association.

The standard GEE approach requires that we have a model for the expected value of the observations given the covariates. With MAR, this marginal model for the mean response will generally not hold for the observed data, so the validity of the analysis is compromised. Methods have been devised for making adjustments to the analysis by using a weighted GEE estimator. The weights have to be estimated using a model for  $\Pr(R_i|Y_i, X_i)$ , hence the non-response model must be explicitly specified and estimated, although the distribution of the error terms need not be. These weighting methods are reviewed in Section 17.5 and discussed in greater detail in Chapter 18.

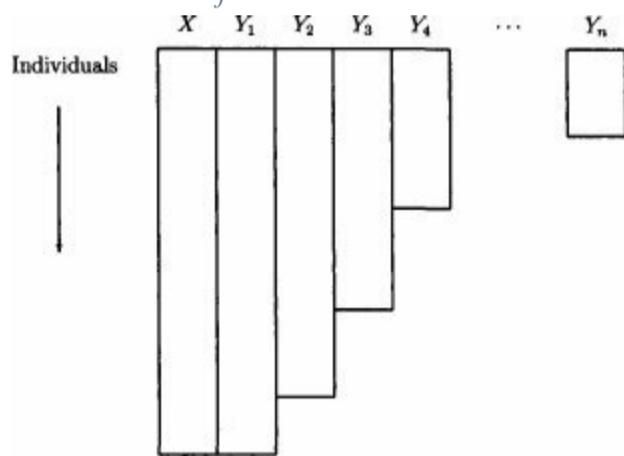
Finally, when longitudinal responses are NMAR, almost all standard methods of longitudinal analysis are not valid. Both GEE methods and standard likelihood-based methods (that ignore the missing data mechanism) yield biased estimates of mean response trends. To obtain valid estimates, joint models for the response and the missing data mechanism are required. Indeed, the term *nonignorable* is used to emphasize that the missing data mechanism must be correctly specified (i.e., cannot be ignored) for inferences about the complete responses. We must also stress that any assumptions made about non-response being NMAR are completely unverifiable from the data at hand. That is, without external (or auxiliary) information about the reasons for missingness or the

missingness mechanism, the observed data provide no information that can either support or refute one NMAR mechanism over another. So, short of tracking down the missing data, any assumptions made about the missingness process are not verifiable. Therefore, when missingness is thought to be NMAR, it is important to carefully assess the sensitivity of inferences to a variety of plausible assumptions concerning the missingness process. However, sensitivity analysis under different assumptions about NMAR missingness is a topic that goes well beyond the scope of this chapter.

## 17.4 DROPOUT

As mentioned earlier, most longitudinal studies are designed to collect data on every individual in the sample at a planned sequence of occasions. However, longitudinal studies habitually suffer from the problem of attrition; that is, some individuals “drop out” of the study prematurely. The term *dropout* refers to the special case where if  $Y_{ik}$  is missing, then  $Y_{ik+1}, \dots, Y_{in}$  are also missing. Alternatively, when expressed in terms of the response indicators, dropout refers to the case where if  $R_{ik} = 0$  then  $R_{ik+1} = \dots = R_{in} = 0$ . This gives rise to the monotone missing data pattern displayed in [Figure 17.1](#), in contrast to the non-monotone patterns that can arise when data are missing intermittently. Note that intermittent missing data give rise to a considerably larger number of potential missing data patterns but, apart from that, do not raise any further technical considerations. As a result the focus of the remainder of this chapter is on dropout.

**Fig. 17.1** Schematic representation of a monotone missing data pattern for dropout, with  $Y_j$  more observed than  $Y_{j+1}$  for  $j = 1, \dots, n - 1$ . Each row represents an individual; the bars represent the subset of individuals with observations on  $Y_j$ .



When there is dropout in a longitudinal study, the key issue is whether those who “drop out” and those who remain in the study differ in any further relevant way. If they do not, then analyses restricted to those remaining in the study yield valid, albeit inefficient, inferences. If they do differ, then such “complete-case” analyses are potentially biased.

In the previous section three different types of missing data mechanisms were distinguished. The same taxonomy can be applied to dropout. That is, dropout can be *completely at random*, *at random*, or *not at random*. When dropout is completely at random the probability of dropout at each occasion is independent of all past, current, and future outcomes (given the covariates). With completely random dropout, an individual leaves the study by a process unrelated to that individual’s outcomes. In contrast, when dropout is at random, the probability of dropout at each occasion can depend on the previously observed outcomes up to, but not including, the current occasion. However, given the observed outcomes, dropout is assumed to be independent of the current and future outcomes. That is, with random dropout the process can depend on the outcomes that have been observed in the past, but given this information, it is unrelated to all future (unobserved) values of the outcome variable following dropout. Finally, when dropout is not at random, the probability of dropping out at each occasion can depend on current and future unobserved outcomes. That is, dropout is said to be not at random when the process depends on the unrecorded values of the outcome variable that would have been observed had the individual remained in the study. In the context of dropout in a longitudinal study, the term “informative” dropout often is used to refer to dropout that is NMAR (similarly “non-informative” dropout often is used to refer to dropout that is either random or completely random). Here the fact of dropout is informative about the distribution of future observations. For example, consider two subjects with the same past history of responses (and covariates) up to time  $t$ . One drops out and the other does not. With MAR, their future observations have the same distribution. In contrast, dropout that is NMAR informs us that the distributions of the future observations will differ. In the general case, nothing in the data can be used to determine the distribution of the future observations of the dropouts; hence the analysis depends strongly on the specification of  $\Pr(R_i|Y_i, X_i)$ .

# Illustration

To emphasize the main distinctions between the three types of dropout mechanism, and their potential impact on a longitudinal analysis, we consider the following simple illustration. Suppose that repeated measurements,  $Y_{it}$  ( $i = 1, \dots, N$ ;  $t = 1, \dots, 5$ ), are generated from a multivariate normal distribution with mean response

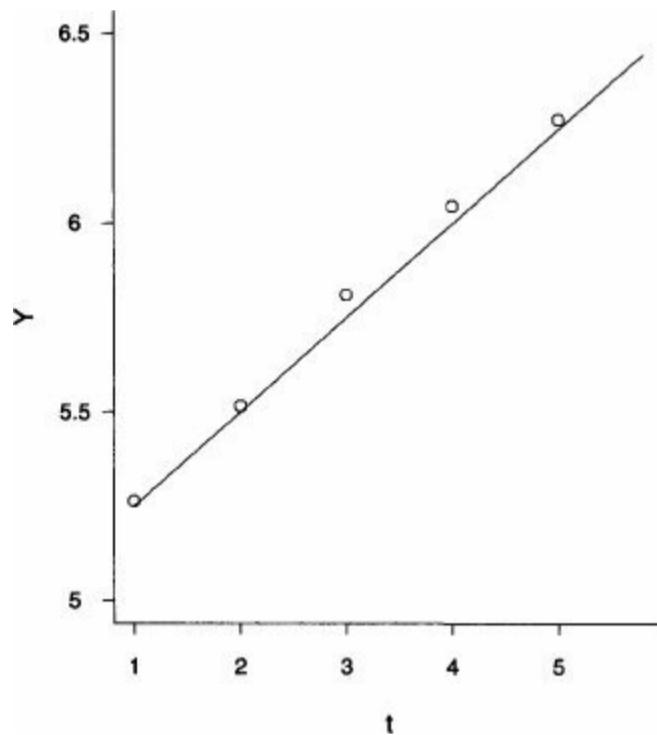
$$E(Y_{it}) = \mu_{it} = \beta_1 + \beta_2 t$$

and covariance

$$\text{Cov}(Y_{is}, Y_{it}) = \rho^{|s-t|}, \text{ for } \rho \geq 0.$$

That is, the variance at each occasion is 1 and assumed to be constant over time, while the correlations have a first-order autoregressive pattern (see Section 7.4). [Figure 17.2](#) displays sample means of simulated data from this model, with  $N = 1000$ ,  $\beta_1 = 5$ ,  $\beta_2 = 0.25$ , and  $\rho = 0.7$ . The sample means show a clear increasing trend over time and virtually coincide with the population regression line (the solid line in [Figure 17.2](#)).

**Fig. 17.2** Population regression line and empirical means at each occasion for simulated complete data.



Next suppose that there is dropout. When there is dropout, we can replace the vector of response indicators,  $R_{it}$  ( $t = 1, \dots, 5$ ), with a simple dropout indicator variable,  $D_i$ , for each individual. The random variable  $D_i$  is recorded for all individuals and  $D_i = k$  if an individual drops out between the  $(k-1)^{th}$  and  $k^{th}$  occasion; that is, only the first  $D_i - 1$  responses are observed. Assume that

$$\log \left\{ \frac{\Pr(D_i = k | D_i \geq k, Y_{i1}, \dots, Y_{ik})}{\Pr(D_i > k | D_i \geq k, Y_{i1}, \dots, Y_{ik})} \right\} = \theta_1 + \theta_2 (Y_{ik-1} - \mu_{ik-1}) + \theta_3 (Y_{ik} - \mu_{ik}).$$

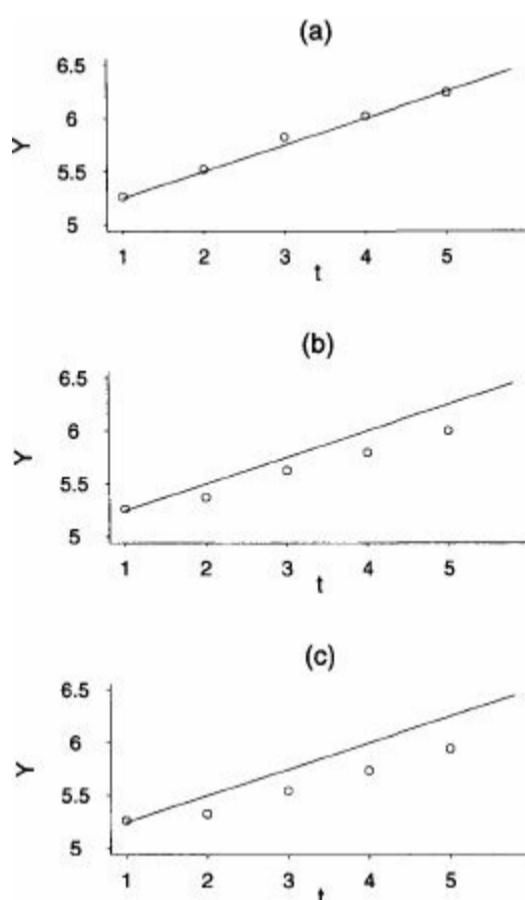
This model specifies that the probability of dropout at any occasion, given dropout has not previously occurred, can depend on the current value and the prior value of the response variable (relative to its mean). We assume that the first response,  $Y_{i1}$ , is fully observed; that is,  $\Pr(D_i = 1 | D_i \geq 1, Y_{i1}) = 0$ . In terms of this model for dropout, consider the following three missing data mechanisms:

- (a) Dropout is MCAR:  $\theta_2 = \theta_3 = 0$ .
- (b) Dropout is MAR:  $\theta_3 = 0$ . and
- (c) Dropout is NMAR:  $\theta_3 \neq 0$ .

[Figure 17.3\(a\)](#) displays simulated data from this model (with  $N = 1000$ ,  $\beta_1 = 5$ ,  $\beta_2 = 0.25$ ,  $\rho = 0.7$ ) when dropout is MCAR (with  $\theta_1 = -0.5$ , and  $\theta_2 = \theta_3 = 0$ ). The conditional probability of dropout at the second through fifth occasions is 0.38 (or  $\frac{e^{-0.5}}{1+e^{-0.5}}$ ). This results in approximately 38% of the responses being missing at the second occasion, 61% missing at the third occasion, 76% missing at the fourth occasion, and 85% missing at the fifth occasion. Despite the large proportion of missing

data, the empirical means at each occasion show a clear linearly increasing trend over time and almost coincide with the population regression line (solid line). Recall that when the missing data mechanism is MCAR, individuals with missing data are a random subset of the sample, and no bias will arise with almost any method of analysis of the observed data (either the complete data or the available data). That is, all of the methods that have been discussed so far will yield valid inferences when missing data are MCAR. To reinforce this point, the estimates of the regression parameters obtained using maximum likelihood (ML) estimation, with correctly specified covariance structure, and using a “working independence” GEE estimator are displayed at the top of [Table 17.2](#). Recall that for a linear model with a “working independence” assumption for the covariance (and a single dispersion parameter), the GEE estimator is identical to the ordinary least squares (OLS) estimator. As expected, both the ML and OLS (or “working independence” GEE) estimates of the intercept and slope are very close to the true values of the population parameters used to generate the data. The minor differences are simply due to sampling variability.

**Fig. 17.3** Population regression line and observed data means at each occasion for simulated data when dropout is (a) completely at random (MCAR), (b) at random (MAR), and (c) not at random (NMAR).



**Table 17.2** Parameter estimates and standard errors for correctly specified likelihood analysis (ML) and “working independence” analysis (OLS/GEE) based on simulated data when dropout is (a) completely at random, (b) at random, and (c) not at random. The true regression parameters are  $\beta_1 = 5.0$  and  $\beta_2 = 0.25$ .

Dropout	Parameter	ML		OLS/GEE	
		Estimate	SE	Estimate	SE <sup>a</sup>
MCAR	Intercept	5.015	0.031	5.022	0.032
	t	0.257	0.016	0.253	0.018
MAR	Intercept	5.003	0.041	5.062	0.043
	t	0.261	0.016	0.182	0.018
NMAR	Intercept	5.058	0.040	5.071	0.043
	t	0.201	0.016	0.162	0.018

<sup>a</sup> Standard errors for OLS/GEE are based on sandwich variance estimator.

[Figure 17.3\(b\)](#) displays simulated data from the same model (with  $N = 1000$ ,  $\beta_1 = 5$ ,  $\beta_2 = 0.25$ ,  $\rho =$

0.7) when dropout is MAR (with  $\theta_1 = -0.5$ ,  $\theta_2 = 0.5$ , and  $\theta_3 = 0$ ). Here dropout at any occasion depends on the previous response but not the current response. Because those with large values (i.e., with previous response  $Y_{ik-1} > \mu_{ik-1}$ ) are more likely to drop out ( $\theta_2 > 0$ ), the empirical means at the second through fifth occasions are discernibly lower than the population regression line. As a result available-data methods such as the GEE will yield biased estimates of mean response trends. In contrast, likelihood-based methods will yield valid estimates when missing data are MAR (or MCAR) and the model for the covariance has been correctly specified. The ML and GEE estimates of the intercept and slope are displayed in the middle of [Table 17.2](#). The ML estimates are very close to the population parameters and only differ due to sampling variability. On the other hand, the “working independence” GEE (or OLS) estimate of the slope shows very discernible bias and underestimates the rate of change over time ( $\hat{\beta}_2 = 0.18$  versus  $\beta_2 = 0.25$ ).

Finally, [Figure 17.3\(c\)](#) displays simulated data from the same model (with  $N = 1000$ ,  $\beta_1 = 5$ ,  $\beta_2 = 0.25$ ,  $p = 0.7$ ) when dropout is NMAR (with  $\theta_1 = -0.5$ ,  $\theta_2 = 0$ , and  $\theta_3 = 0.5$ ). Here dropout at any occasion depends on the current value of the response. Because those with large values at a given occasion are more likely to be unobserved at that occasion ( $\theta_3 > 0$ ), the empirical means are discernibly lower than the population regression line. As a result available-data methods such as the GEE yield biased estimates of mean response trends. Furthermore likelihood-based methods that ignore the missing data mechanism also yield biased estimates of mean response trends. To reinforce this point, the ML and GEE estimates of the regression parameters are displayed at the bottom of [Table 17.2](#). The ML and GEE estimates of the slope show large biases. Of note, the magnitude of the bias is somewhat smaller for ML; however, this cannot be expected in general unless the correlation among the responses is very high. When the correlation among the responses is very high and dropout at any occasion depends only on the current value of the response, the dropout mechanism can often be approximated by an ignorable dropout mechanism that conditions on all previously observed responses

$$\Pr(D_i = k | D_i \geq k, Y_{ik}) \approx \Pr(D_i = k | D_i \geq k, Y_{i1}, \dots, Y_{ik-1}).$$

For example, when the data are simulated from the same model but with a correlation parameter  $\rho = 0.9$  instead of  $\rho = 0.7$ , the ML estimate of the slope is  $\hat{\beta}_2 = 0.24$  and the magnitude of the bias is significantly reduced. In contrast, the OLS/GEE estimate of the slope is  $\hat{\beta}_2 = 0.11$  and remains highly biased under this NMAR dropout mechanism.

## 17.5 COMMON APPROACHES FOR HANDLING DROPOUT

In this section we present a short review of some of the most commonly used methods for handling dropout in longitudinal analysis. We also discuss the assumptions about dropout required for each of the methods to yield valid inferences. We note that many traditional methods for handling missing data (e.g., complete-case analysis, imputation) became popular when the only approaches for analyzing data were ones based on complete and balanced data.

# Complete-Case Analysis

One approach to handling dropout is to simply exclude all data from the analysis on any subject who drops out. That is, a so-called complete-case analysis can be performed by excluding any subjects that do not have data at all intended measurement occasions. We must stress that this method is very problematic and is rarely an acceptable approach to the analysis. It will yield unbiased estimates of mean response trends only when it can be assumed that dropout is MCAR. Recall that when dropout is MCAR, the study “completers” are a random subsample of the original sample from the population. However, even in cases where the MCAR assumption might be tenable, a complete-case analysis is very unappealing because of the reduction in the number of subjects contributing to the analysis. A complete-case analysis can be immensely inefficient, leading to an analysis with reduced statistical power.

## Available-Data Analysis

Another approach for handling dropout is the available-data method. This is not a single method, but a very general term that refers to a wide collection of techniques that can readily incorporate vectors of repeated measures of unequal length in the analysis. For example, standard applications of generalized least squares (GLS) or the generalized estimating equations approach can be considered available-data methods, since these approaches base the analysis on all of the available observations. In general, available-data methods are more efficient than complete-case methods because they incorporate the partial information obtained from those who drop out. However, many available-data methods will yield valid analyses only if the conditional (i.e., conditional on  $X_i$ ) means and covariances of the observed components of  $Y_i$  among those who drop out coincide with the corresponding conditional means and covariances of  $Y_i$  in the target population. As a result available-data methods will yield biased estimates of mean response trends unless dropout is MCAR. In general, for available-data (and complete-case) methods to be valid we require that dropout is MCAR.

# Imputation

A third approach, and one that is widely used in practice, is some form of imputation for the missing responses following dropout. The idea behind imputation is very simple: substitute or fill in the values that were not recorded with imputed values. One of the chief attractions of imputation methods is that, once a filled-in data set has been constructed, standard methods for complete data can be applied. However, methods that rely on just a single imputation, creating only a single filled-in data set, fail to acknowledge the uncertainty inherent in the imputation of missing data. Multiple imputation circumvents this difficulty. In multiple imputation the missing values are replaced by a set of  $m$  plausible values, thereby acknowledging the uncertainty about what values to impute for the missing responses. The  $m$  filled-in data sets produce  $m$  different sets of parameter estimates and their standard errors. These are then appropriately combined to provide a single estimate of the parameters of interest, together with standard errors that reflect the uncertainty inherent in the imputation of the unobserved responses. Typically a small number of imputations, for instance,  $10 \leq m \leq 25$ , is sufficient to obtain realistic estimates of the sampling variability. A more detailed description of multiple imputation is given in Chapter 18.

Although the main idea behind imputation is very simple, what is less clear-cut is how to produce the imputed values for the missing responses. Next we consider some of the commonly used methods for imputing missing data. One widely used imputation method, especially in longitudinal clinical trials, is “last value carried forward” (LVCF), occasionally referred to as “last observation carried forward” (LOCF). This is a single imputation method that fills-in or imputes the missing values following dropout with the last observed value for that subject. Despite its widespread use, it should be recognized that LVCF makes a strong, and often very unrealistic, assumption that the responses following dropout remain constant at the last observed value prior to dropout. Perhaps the only setting where this assumption might conceivably be appropriate is when dropout is due to recovery or cure. In the context of placebo-controlled longitudinal clinical trials, there appears to be some statistical folklore that LVCF yields a *conservative* estimate of the comparison of an active treatment versus the control. However, this is a gross misconception, and will only be true to the extent that the active treatment prior to dropout has carry-over effects following dropout. In many clinical trials this is unlikely to be the case; instead, dropout from the active treatment (e.g., due to adverse side effects) might very well result in a deterioration of the response.

Despite frequent and well-founded criticisms by statisticians, LVCF is still widely used to handle dropouts in clinical trials. Regulatory agencies such as the U.S. Food and Drug Administration (FDA) seem to encourage the continuing use of LVCF as a method for handling dropouts, despite all of its obvious shortcomings. Except in very rare cases (as mentioned above), we do not recommend the use of LVCF as a method for handling dropout. In Section 17.6 we provide an illustration of the bias that can arise when LVCF is used to impute missing values, highlighting that LVCF does not necessarily yield a *conservative* estimate of the comparison of an active treatment versus the control; for readers who find the level of detail in this section challenging, Section 17.6 can be omitted at first reading without loss of continuity.

Variations on the LVCF theme include baseline value carried forward and worst value carried forward. Worst value carried forward is most often used in comparisons of an active treatment to a placebo, since it is assumed to be conservative in that setting. However, both of these alternatives suffer the same difficulties as LVCF and cannot be counted on to give unbiased treatment estimates. In addition all of the methods suffer from optimistic standard error estimates. It is easy to see that these analyses give smaller standard errors than complete-case, or even available-data estimates because they assume complete data on everyone. However, they will generally give smaller standard errors than what we would expect if we had been fortunate enough to have complete data on everyone. This is because the variability of baseline measurements is usually smaller because of selection criteria into the study, and as we move out in time, the observations tend to become more variable. Hence substituting baseline or intermediate values for final values can be expected to give a less variable data set. It is also true if we use worst value, since worst values are often similar

especially for responses based on a scale. Therefore we caution that neither LVCF nor any of its variants, such as baseline value carried forward, provide a legitimate approach for analyzing incomplete data. These ad hoc methods of imputation typically produce bias whose direction and magnitude depend on both the true, but unknown, treatment effect and the dropout rates in the treatment groups (see Section 17.6). In addition, similar to other single imputation methods, these methods artificially increase the amount of information in the data by regarding imputed and actually observed values on an equal footing.

Other imputation methods that have a much firmer theoretical foundation draw values of  $Y_i^M$  from the conditional distribution of the missing responses given the observed responses,  $f(Y_i^M|Y_i^O, X_i)$ . With the monotone missing data patterns produced by dropouts, it is relatively straightforward to impute missing values by drawing values of  $Y_i^M$  from  $f(Y_i^M|Y_i^O, X_i)$  in a sequential manner. A variety of imputation methods can be used to draw values from  $f(Y_i^M|Y_i^O, X_i)$ ; we describe two distinct methods in our discussion of multiple imputation in Chapter 18. When missing values are imputed from  $f(Y_i^M|Y_i^O, X_i)$ , regardless of the particular imputation method adopted, subsequent analyses of the observed and imputed data are valid for dropouts that are MAR (or MCAR). Furthermore multiple imputation ensures that the uncertainty is properly accounted for.

Finally, there is another related form of imputation where the missing responses are effectively imputed by modeling and estimating parameters for the joint distribution of  $Y_i, f(Y_i|X_i)$ . When dropout is MCAR or MAR, likelihood-based methods can be used based solely on the marginal distribution of the observed data. If dropout is MCAR or MAR (and the parameters of the dropout and outcome processes are distinct, a technical requirement that can usually be assumed in practice), then ML estimates can be obtained by maximizing  $f(Y_i^O|X_i)$ , where  $f(Y_i^O|X_i)$  denotes the ordinary marginal distribution of the particular subset of  $Y_i$  determined by  $Y_i^O$ . Importantly, likelihood-based inference does not require specification of the dropout mechanism and the contribution of  $\Pr(D_i|Y_i^O, X_i)$  to the likelihood can be ignored. In a certain sense, the missing values are validly predicted by the observed data via the model for the conditional mean,  $E(Y_i^M|Y_i^O, X_i)$ . In many missing data situations, ML estimates of the parameters can be easily obtained by an iterative *EM algorithm* that alternates between filling in missing values (the expectation or E-step), then maximizing the likelihood for the resulting filled-in data set (the maximization or M-step). For example, if the responses are assumed to have a multivariate normal distribution, then predictions of the missing values are based on the conditional mean of  $Y_i^M$ , given  $Y_i^O$  (and  $X_i$ ),

$$E(Y_i^M|Y_i^O, X_i) = \mu_i^M + \Sigma_i^{MO} \Sigma_i^{O^{-1}} (Y_i^O - \mu_i^O),$$

where  $\mu_i^M$  and  $\mu_i^O$  denote those components of the mean response vector corresponding to  $Y_i^M$  and  $Y_i^O$ , and  $\Sigma^O$  and  $\Sigma_i^{MO}$  denote those components of the covariance matrix corresponding to the covariance among the elements of  $Y_i^O$  and the covariance between  $Y_i^M$  and  $Y_i^O$ . This simple implementation of the EM algorithm works for estimating means in the setting of the multivariate normal distribution, but for estimating variances or covariances, somewhat more complex expressions are required for filling in the missing observations.

# Weighting Methods

An alternative approach for handling dropout is to weight the observed data in some appropriate way. In weighting methods, the under-representation of certain response profiles in the observed data is taken into account and corrected. A variety of different weighting methods that adjust for dropout have been proposed. These approaches are often called propensity weighted or inverse probability weighted (IPW) methods. Here the underlying idea is to base estimation on the observed responses but weight them to account for the probability of remaining in the study. The probability of remaining in the study can be estimated as a function of the observed responses prior to dropout, and also as a function of the covariates and any additional variables or subject characteristics that are thought likely to predict dropout.

In the simplest version of this approach a single weight,  $W_i$ , is calculated only for those individuals who complete the study. The weight for each individual denotes the inverse probability of remaining in the study until the last intended measurement occasion; that is,  $w_i = \{\Pr(D_i = n + 1)\}^{-1}$ . It can be computed sequentially as the inverse of the following product of the conditional probabilities of remaining in the study at each occasion:

$$\begin{aligned} w_i &= \{\Pr(D_i = n + 1)\}^{-1} \\ &= \{\Pr(D_i > 1|D_i \geq 1) \times \Pr(D_i > 2|D_i \geq 2) \times \cdots \times \Pr(D_i > n|D_i \geq n)\}^{-1} \\ &= (\pi_{i1} \times \pi_{i2} \times \cdots \times \pi_{in})^{-1}, \end{aligned}$$

where  $\pi_{ik} = \Pr(D_i > k|D_i \geq k)$  can be estimated from those remaining at the  $(k - 1)^{th}$  occasion, given the recorded history of all available data up to the  $(k - 1)^{th}$  occasion. Given the estimated weight,  $\hat{w}_i$ , a weighted complete-case analysis can be performed. For example, the GEE approach can be adapted to handle data that are MAR by making adjustments to the analysis for the probability of remaining in the study. One variant of this approach is to use a weighted GEE to analyze the data from the study “completers”, with weights inversely proportional to the estimated probability that the  $i^{th}$  subject completes the study. In the weighted complete-case analysis each subject’s contribution to the analysis is weighted by  $\hat{w}_i$ , thereby providing valid estimates of the mean response trends when dropout is MAR.

Inverse probability weighted methods were first proposed in the sample survey literature where the weights are known and based on the survey design. In the sample survey setting, units are sampled with unequal probability of selection and therefore must be given correspondingly unequal weights (inverse probability weighting) in the analysis. In the missing data setting, the intuition behind the weighting methods is that each subject’s contribution to the weighted complete-case analysis is replicated  $w_i$  times, in order to count once for herself and  $(w_i - 1)$  times for those subjects with the same history of responses and covariates who do not complete the study. These weights correct for the under-representation of certain response profiles in the observed data due to dropout. The weighting methods are valid provided that the model that produces the estimated  $w_i$  is correctly specified.

In longitudinal analyses,  $w_i$  is not ordinarily known, but must be estimated from the observed data (e.g., using a repeated sequence of logistic regressions for the  $\pi_{ik}$ ’s). Therefore the variance of inverse probability weighted estimators should also account for estimation of  $w_i$ . Counter-intuitively, estimation of the weights from the data at hand leads to improvements in precision. Finally, we note that this approach for handling dropout can be made more efficient by conducting an appropriately weighted available-data analysis. This requires that occasion-specific weights for each individual,  $w_{ij}$ , be incorporated into the analysis, where  $w_{ij}$  denotes the inverse probability that the  $i^{th}$  subject is still in the study at the  $j^{th}$  occasion. A more detailed description of inverse probability weighting methods is given in Chapter 18.

# 17.6 BIAS OF LAST VALUE CARRIED FORWARD IMPUTATION\*

In Section 17.5 we noted that last value carried forward (LVCF) imputation makes a strong, and often very unrealistic, assumption that the responses following dropout remain constant at the last observed value prior to dropout. In this section<sup>1</sup> we provide a demonstration of the bias that can arise when LVCF is used to impute missing values, highlighting that LVCF does not necessarily yield a *conservative* estimate of the comparison of an active treatment to the control. When LVCF is used in the comparison of an active treatment to a control, the bias can go in either direction. The content of this section is somewhat technical and can be omitted on first reading of this chapter without loss of continuity.

The three mechanisms, MCAR, MAR, and NMAR, introduced in Section 17.2 are only one approach to modeling missing data, albeit the most widely used approach. There is another approach known as *pattern mixture models* that can be easier to understand and specify in certain missing data situations, especially when modeling dropout. Simply put, pattern mixture models stratify the data by missing data patterns and assume a different model for the data within each stratum (e.g., assume that subjects who drop out are those who are likely not responding well to treatment). To study the potential bias of LVCF imputation, it is easier to use a pattern mixture model approach to determine what kind of bias might accrue.

To illustrate the potential bias that can arise from LVCF, consider a simple two-group design (e.g., active treatment versus control) with two repeated measures of the response, one at baseline, the other at end of follow-up. Let  $\text{trt}_i = 1$  if the  $i^{\text{th}}$  subject is assigned to the active treatment group and  $\text{trt}_i = 0$  if assigned to the control group. We assume that the baseline response,  $Y_{i1}$ , is always observed ( $R_{i1} = 1$  for all individuals) and the probability of the follow-up response,  $Y_{i2}$ , being observed is  $\pi_0 = \Pr(R_{i2} = 1|\text{trt}_i = 0)$  for those in the control group and  $\pi_1 = \Pr(R_{i2} = 1|\text{trt}_i = 1)$  for those in the active treatment group. When stratified in terms of being a “dropout” (with  $R_{i2} = 0$ ) or a “completer” (with  $R_{i2} = 1$ ), two saturated (or unrestricted) models for the change in the mean response can be specified as

$$(17.4) E(Y_{i2} - Y_{i1}|\text{trt}_i, R_{i2} = 0) = \alpha_1 + \alpha_2 \text{trt}_i,$$

$$(17.5) E(Y_{i2} - Y_{i1}|\text{trt}_i, R_{i2} = 1) = \gamma_1 + \gamma_2 \text{trt}_i.$$

The alert reader would have recognized that without some assumptions about the missing data mechanism,  $\alpha_1$  and  $\alpha_2$  in (17.4) cannot be estimated from the available data because  $Y_{i2}$  is not observed for those who drop out. Recall from Section 17.2 that for the special case where the missing data mechanism is assumed to be MCAR, the distribution of  $Y_i$  (given the covariates) is the same in each of the distinct sub-populations defined by the missing data patterns. In this illustration there are only two sub-populations, the “dropouts” and the “completers.” Therefore, under the MCAR assumption,  $\gamma_k = \alpha_k$  for  $k = 1, 2$  in equations (17.4) and (17.5).

In such a study design, the primary goal of the analysis is to compare the two groups in terms of their changes in response from baseline to follow-up. Specifically, the parameter of primary interest can be expressed as

$$(17.6) \delta = E(Y_{i2} - Y_{i1}|\text{trt}_i = 1) - E(Y_{i2} - Y_{i1}|\text{trt}_i = 0).$$

Note that the parameter  $\delta$  is expressed in terms of  $E(Y_{i2} - Y_{i1}|\text{trt}_i)$  not  $E(Y_{i2} - Y_{i1}|\text{trt}_i, R_{i2})$ . We can obtain  $E(Y_{i2} - Y_{i1}|\text{trt}_i)$  by simply averaging  $E(Y_{i2} - Y_{i1}|\text{trt}_i, R_{i2})$  over the distribution of  $R_{i2}$  (given treatment group). Specifically,

$$\begin{aligned} E(Y_{i2} - Y_{i1}|\text{trt}_i) &= E(Y_{i2} - Y_{i1}|\text{trt}_i, R_{i2} = 0) \Pr(R_{i2} = 0|\text{trt}_i) \\ &\quad + E(Y_{i2} - Y_{i1}|\text{trt}_i, R_{i2} = 1) \Pr(R_{i2} = 1|\text{trt}_i). \end{aligned}$$

From equations (17.4) and (17.5), it can be seen that

$$\begin{aligned}
E(Y_{i2} - Y_{i1} | \text{trt}_i = 1) &= E(Y_{i2} - Y_{i1} | \text{trt}_i = 1, R_{i2} = 0) \Pr(R_{i2} = 0 | \text{trt}_i = 1) \\
&\quad + E(Y_{i2} - Y_{i1} | \text{trt}_i = 1, R_{i2} = 1) \Pr(R_{i2} = 1 | \text{trt}_i = 1) \\
&= E(Y_{i2} - Y_{i1} | \text{trt}_i = 1, R_{i2} = 0)(1 - \pi_1) \\
&\quad + E(Y_{i2} - Y_{i1} | \text{trt}_i = 1, R_{i2} = 1)\pi_1 \\
&= (\alpha_1 + \alpha_2)(1 - \pi_1) + (\gamma_1 + \gamma_2)\pi_1.
\end{aligned}$$

Similarly

$$\begin{aligned}
E(Y_{i2} - Y_{i1} | \text{trt}_i = 0) &= E(Y_{i2} - Y_{i1} | \text{trt}_i = 0, R_{i2} = 0) \Pr(R_{i2} = 0 | \text{trt}_i = 0) \\
&\quad + E(Y_{i2} - Y_{i1} | \text{trt}_i = 0, R_{i2} = 1) \Pr(R_{i2} = 1 | \text{trt}_i = 0) \\
&= E(Y_{i2} - Y_{i1} | \text{trt}_i = 0, R_{i2} = 0)(1 - \pi_0) \\
&\quad + E(Y_{i2} - Y_{i1} | \text{trt}_i = 0, R_{i2} = 1)\pi_0 \\
&= \alpha_1(1 - \pi_0) + \gamma_1\pi_0.
\end{aligned}$$

Therefore

$$\begin{aligned}
\delta &= \{(\alpha_1 + \alpha_2)(1 - \pi_1) + (\gamma_1 + \gamma_2)\pi_1\} \\
&\quad - \{\alpha_1(1 - \pi_0) + \gamma_1\pi_0\} \\
&= \alpha_2 + (\pi_1 - \pi_0)(\gamma_1 - \alpha_1) + \pi_1(\gamma_2 - \alpha_2).
\end{aligned}$$

Next we consider the target parameter of the analysis based on LVCF imputation. Recall that in the LVCF analysis it is assumed that  $E(Y_{i2} | \text{trt}_i, R_{i2} = 0) = E(Y_{i1} | \text{trt}_i, R_{i2} = 0)$ ; that is,  $E(Y_{i2} - Y_{i1} | \text{trt}_i, R_{i2} = 0) = 0$ . Under this assumption it can be shown that the target parameter for the LVCF analysis, denoted  $\delta_{LVCF}$ , is

$$\delta_{LVCF} = (\gamma_1 + \gamma_2)\pi_1 - \gamma_1\pi_0 = (\pi_1 - \pi_0)\gamma_1 + \pi_1\gamma_2.$$

Thus it is transparent that  $\delta_{LVCF} \neq \delta$  and that, in general, the LVCF analysis will yield a biased estimate of  $\delta$ , the parameter of primary interest.

To see what kind of bias might accrue, consider the case where the underlying missing data mechanism is assumed to be MCAR. Under the MCAR assumption,  $\gamma_k = \alpha_k$  for  $k = 1, 2$  in [equations \(17.4\)](#) and [\(17.5\)](#), and the expression for  $\delta$  simplifies to  $\delta = \gamma_2 = \alpha_2$ . However, even under the MCAR assumption, the target parameter of the LVCF analysis is

$$\delta_{LVCF} = (\pi_1 - \pi_0)\gamma_1 + \pi_1\gamma_2 \neq \delta.$$

Thus the LVCF analysis yields a biased estimate of  $\delta$  even under MCAR. Under MCAR the amount and direction of bias for the LVCF analysis is given by

$$(\delta_{LVCF} - \delta) = (\pi_1 - \pi_0)\gamma_1 - (1 - \pi_1)\gamma_2.$$

To demonstrate that the bias can be in any direction, both positive and negative, consider the case where  $\gamma_1 \neq 0$  and  $\gamma_2 = 0$ . This corresponds to the scenario where there is change in the mean response from baseline in both treatment groups but at the same rate (i.e., there is no differential treatment effect on the pattern of change). Then the bias can go in either direction depending on the signs of  $(\pi_1 - \pi_0)$  and  $\gamma_1$ . This highlights how the LVCF analysis can potentially produce an apparent treatment effect, favoring either the active treatment group or the control group, when no such effect exists ( $\delta = \gamma_2 = 0$ ). Under the assumption that the missing data are MAR, it is also possible to derive expressions for the bias of the LVCF analysis. Under MAR, expressions for the bias are somewhat more complicated but nonetheless reveal that the LVCF analysis also yields a biased treatment comparison, with bias that can operate in either direction.

In summary, although LVCF is a widely used imputation method, especially in longitudinal clinical trials, it makes a strong, and often very unrealistic, assumption about the responses following dropout. As we have seen, even when missingness can be assumed to be MCAR, LVCF yields a biased treatment comparison and the bias can go in either direction. Moreover, due to this bias, an LVCF analysis can potentially yield an apparent treatment effect when no such effect exists. In contrast, under MCAR, almost any other method of analysis, including a complete-case analysis, is unbiased. Therefore, except in very rare cases (e.g., when dropout is due to cure or recovery) we do not recommend the use of LVCF as a method for handling dropout. Finally, we note that LVCF happens to be equivalent to “baseline value carried forward” in the simple illustration used in this section. Therefore, all our criticism of LVCF applies equally to “baseline value carried forward” imputation.

## 17.7 FURTHER READING

A useful discussion of methods for handling dropout in longitudinal studies can be found in Heyting et al. (1992). The tutorial article by Hogan et al. (2004) provides a comprehensive overview of more recent developments in methods for adjusting for drop-out within likelihood-based and semiparametric modeling frameworks and illustrates their application with two worked examples. White et al. (2011) suggest a general framework for “intention to treat” analysis in randomized clinical trials that depends on making plausible assumptions about the missing data and including all participants in sensitivity analyses.

In longitudinal clinical trials, “last value carried forward” (LVCF) imputations are still widely used to handle dropouts; see Ware (2003), Cook et al. (2004), Molenberghs et al. (2004), and Kenward and Molenberghs (2009) for critiques of this method and its variants (e.g., baseline value carried forward). The illustration of the potential bias that can arise from LVCF in Section 17.6 is based on similar derivations in Chapter 27 of Molenberghs and Verbeke (2005).

# Bibliographic Notes

Rubin (1976) developed the taxonomy for describing the assumptions concerning the dependence of the missingness process on observed and unobserved responses. Little and Rubin (2001) is the definitive textbook on missing data, providing a comprehensive description of the theory and application of methods for handling missing data; also see Schafer (1997), Tsiatis (2006), Molenberghs and Kenward (2007), and Daniels and Hogan (2008). Laird (1988) discusses missing data issues in longitudinal studies; also see the review articles by Little (1995) and Kenward and Molenberghs (1999). The EM algorithm, a general technique for ML estimation with incomplete data, is discussed in the seminal paper by Dempster et al. (1977). Finally, Hogan and Laird (1996) and Little and Yau (1996) discuss methods for handling missing data for “intention to treat” analysis in randomized clinical trials.

<sup>1</sup> This section derives formulae for the potential bias of LVCF in a simple setting and is based on similar derivations in Chapter 27 of Molenberghs and Verbeke (2005).

# *Chapter 18*

## *Missing Data and Dropout: Multiple Imputation and Weighting Methods*

### **18.1 INTRODUCTION**

In the previous chapter we distinguished between different types of missing data mechanisms by their assumptions concerning whether missingness is related to observed and unobserved responses. We emphasized that conventional likelihood-based analyses of incomplete data (e.g., linear mixed effects models fitted using standard statistical software such as PROC MIXED in SAS, the `lme` function in the `nlme` package in R and S-Plus, and the `xtmixed` command in Stata) yield valid inferences when data are MAR if the joint distribution of the vector of responses has been correctly specified. When the joint distribution is assumed to be multivariate normal, this requires correct specification of not only the model for the mean response, but also the model for the covariance among the responses. In this chapter we discuss two alternative approaches for handling missing data: multiple imputation and weighting methods. Both approaches are appealing in settings where a conventional likelihood-based analysis is no longer straightforward. For example, when there are missing covariates as well as missing responses, a likelihood-based analysis is not straightforward and cannot be implemented using standard statistical software. Also there is no convenient specification of the joint distribution of the vector of responses for marginal models when the responses are discrete; instead, the generalized estimating equations (GEE) approach is routinely used and requires the stronger assumption of MCAR. In both of these settings, as well as others, multiple imputation and weighting methods provide a convenient and practical method for handling missing data in longitudinal analyses.

## 18.2 MULTIPLE IMPUTATION

Multiple imputation is a flexible method for handling missing data that has recently been implemented in numerous commercially available software packages (e.g., SAS, Stata, SPSS, R, and S-Plus), as well as in more specialized software programs (e.g., SUDAAN and Solas). As discussed in Chapter 17, imputation is an intuitively simple technique: we “fill in” or impute plausible values for the missing data, thereby creating a “completed” data set that can be analyzed using standard statistical methods for complete data. Imputation allows us to proceed with the analysis of the completed data set as though there were no missing data at all. However, if each missing datum is filled in with one plausible value only, any subsequent analysis of this single completed data set is problematic. The trouble with such an analysis is that the imputed values are implicitly treated as though they are known, neglecting the fact that there is inherent uncertainty surrounding the imputed values. Conventional methods for standard error estimation do not properly account for this uncertainty. Specifically, any analysis of the single imputed data set as though it were a complete data set will, in general, produce anti-conservative results. By “anti-conservative,” we mean that nominal  $p$ -values will tend to be too small and confidence intervals will be too narrow. Multiple imputation was developed to correct this problem.

With multiple imputation methods each missing datum is filled in with plausible values multiple times, producing multiple completed data sets. The replacement of missing data with multiple plausible values ensures that the uncertainty associated with the imputed values can be properly accounted for. Typically the number of imputations is relatively small, say between 10 and 25. Each of the completed data sets is then analyzed using standard methods for complete data, as if there were no missing data. These analyses produce a set of results that are then appropriately combined to provide a single estimate of the parameters of interest, together with standard errors that reflect the uncertainty inherent in the imputation of the missing data. Specifically, if we assume that  $m > 1$  filled-in data sets are created, then  $m$  different estimates of the regression parameters  $\beta$ , say  $\hat{\beta}^{(k)}$  (for  $k = 1, \dots, m$ ), can be obtained from the separate analyses of each of the  $m$  data sets. In addition the  $m$  analyses of the filled-in data sets also yield  $m$  estimates of the covariance of  $\hat{\beta}^{(k)}$ , for  $k = 1, \dots, m$ . The multiple imputation estimate of  $\beta$  is simply the unweighted average of the  $m$  estimates,

$$\hat{\beta} = \bar{\beta} = \frac{1}{m} \sum_{k=1}^m \hat{\beta}^{(k)}.$$

The estimated covariance of  $\hat{\beta}$  is given by

$$\widehat{\text{Cov}}(\hat{\beta}) = W + (1 + m^{-1})B,$$

where

$$W = \frac{1}{m} \sum_{k=1}^m \widehat{\text{Cov}}(\hat{\beta}^{(k)})$$

and

$$B = \frac{1}{m-1} \sum_{k=1}^m (\hat{\beta}^{(k)} - \bar{\beta})(\hat{\beta}^{(k)} - \bar{\beta})'.$$

Although the expression for  $\widehat{\text{Cov}}(\hat{\beta})$  appears somewhat complicated, it simply combines two inherent sources of variability: the within-imputation variability ( $W$ ) and the between-imputation variability ( $B$ ).

So far, we have briefly described how to make inferences from a statistical analysis based on multiple imputation. It can be thought of as a three-step process. First, the missing data are filled in  $m$  times to create  $m$  completed data sets. Second, the  $m$  completed data sets are analyzed using standard statistical methods (e.g., fitting any of the regression models discussed in Parts II and III). Finally, the results from the  $m$  analyses of the completed data sets are combined in the manner described above. For the remainder of this section, we focus on the first step and consider how to create multiple completed data sets using methods that incorporate appropriate variability across the  $m$  imputations. At the outset it must be recognized that there are numerous different methods of imputation, although

most have a similar basis. Moreover the method of choice is often determined by the patterns of missingness (e.g., monotone versus non-monotone patterns of missingness) in the data set at hand. Although there are numerous methods of imputation, one general principle should guide the choice of imputation method: “proper” imputations should be drawn at random from the conditional (or so-called predictive) distribution of the missing data given the observed data. In general, a proper imputation of  $Y_i^M$  is obtained by *randomly* drawing values from  $f(Y_i^M|Y_i^O, X_i)$ . By choosing to draw values from  $f(Y_i^M|Y_i^O, X_i)$ , we are implicitly assuming that missingness is MAR; that is, the predictive distribution of the missing data, given the observed data, does not depend on the observed response pattern,  $R_i$ , with  $f(Y_i^M|Y_i^O, X_i, R_i) = f(Y_i^M|Y_i^O, X_i)$ . To randomly sample values from  $f(Y_i^M|Y_i^O, X_i)$ , it is important to distinguish two settings: (1) monotone missing data patterns, and (2) non-monotone missing data patterns. Imputation is far more straightforward in the former case, and much of the remainder of this section focuses on monotone missing data. For the latter case, iterative computational methods are usually required.

Before describing specific methods for imputing longitudinal data with missing responses, it is worth noting that these methods implicitly assume the data set is structured in a “wide” rather than a “long” format, with a single “record” for each individual. In a “wide” format, methods for imputation of a missing response at any particular occasion can exploit the positive correlation with the responses at any of the remaining occasions. This is discussed in greater detail in Section 18.6.

## 18.2.1 Monotone Missing Data Patterns

As noted in the previous chapter, monotone missing data patterns arise in longitudinal studies when missingness occurs only through dropout. For example, suppose that the first response,  $Y_1$ , is always observed but subsequent responses are missing due to dropout. Then a monotone missing pattern is produced where  $Y_1$  is fully observed, the second response,  $Y_2$ , has the fewest missing values,  $Y_3$  has the second fewest missing values, and so on. With monotone missing data patterns, missing values can be imputed by first fitting an appropriate model (e.g., a regression model) to predict  $Y_{i2}$  from  $Y_{i1}$  and  $X_i$  and then randomly sampling from this model to impute the missing values in the second response. Next, based on an appropriate model, the missing values in the third response can be imputed by predicting  $Y_{i3}$  from  $X_i$ ,  $Y_{i1}$ , and both the *observed* and *imputed* values of  $Y_{i2}$ . Imputation of the remaining missing values can continue in a similar way until all of the missing values have been filled in. The resulting set of imputed values is a proper imputation of  $Y_i^M$  from  $f(Y_i^M|Y_i^O, X_i)$  when the MAR assumption holds. There are two commonly used methods of imputation for monotone missing data patterns: regression methods and predictive mean matching methods. We briefly describe each approach in turn.

# Regression Methods

In regression methods for imputing longitudinal data, the missing responses are imputed sequentially using all preceding responses in the monotone pattern (and any subset of the covariates) as “predictors” in a regression model. That is, a series of regression models for  $Y_{ik}$ , given  $Y_{i1}, \dots, Y_{ik-1}$  and  $X_i$  (or a subset of  $X_i$ ), are fitted to the observed data. When the responses are continuous, standard linear regression is widely used to generate imputations. For example, when  $Y_{ik}$  is continuous, a linear regression model

$$E(Y_{ik}|Y_{i1}, \dots, Y_{ik-1}, X_i) = \gamma_1 + \gamma_2 Y_{i1} + \dots + \gamma_k Y_{ik-1},$$

can be fitted using the observed data on subjects who have not dropped out by the  $k^{th}$  occasion; alternative models may be needed when linearity is insufficient to capture the functional forms for the relationships. For simplicity, we have assumed no dependence on  $X_i$  in the regression model above. More generally, a series of linear regression models

$$E(Y_{ik}|Y_{i1}, \dots, Y_{ik-1}, X_i) = Z'_{ik}\gamma,$$

can be fitted to data on the  $N_k$  subjects who have not dropped out by the  $k^{th}$  occasion. In a departure from the notation used in previous chapters, we let  $Z_{ik}$  denote a vector composed of  $Y_{i1}, \dots, Y_{ik-1}$  and any subset of the components of  $X_i$ . In this model  $\gamma$  denotes a  $q \times 1$  vector of regression parameters relating  $Y_{ik}$  to the preceding responses and covariates. Also, in a slight abuse of notation, note that there is a separate set of regression parameters,  $\gamma$ , for the regression model at each occasion and that the corresponding dimensions of  $Z_{ik}$  and  $\gamma$  can vary in the models for different occasions. The linear regression models introduced here, however, are for the purpose of imputation and are not the same models that were discussed in Part II for longitudinal analyses. The regression models considered here are imputation and not analysis models.

For the regression model at the  $k^{th}$  occasion,

$$E(Y_{ik}|Y_{i1}, \dots, Y_{ik-1}, X_i) = Z'_{ik}\gamma,$$

the regression parameters (and the residual variance) can be estimated via ordinary least squares (OLS). The fitted linear regression produces estimates of the regression parameters,  $\hat{\gamma}$ , and their associated covariance matrix,

$$\widehat{\text{Cov}}(\hat{\gamma}) = \hat{\sigma}^2 \left( \sum_{i=1}^{N_k} Z_{ik} Z'_{ik} \right)^{-1},$$

where  $\hat{\sigma}^2$  is an estimate of the residual variance and  $Z_{ik}$  is the design vector for the regression of  $Y_{ik}$  on  $Y_{i1}, \dots, Y_{ik-1}$ , and any subset of the components of  $X_i$ .

In principle, the fitted regression could be used to produce predictions or imputations of the missing values of  $Y_{ik}$ . However, because the fitted regression produces a deterministic prediction of the missing values on  $Y_{ik}$  for any fixed  $Z_{ik}$ , we must incorporate random variation to reflect the uncertainty of the imputations. Specifically, we need to add to the predicted value for  $Y_{ik}$  a random draw from the residual distribution of  $Y_{ik}$  given  $Y_{i1}, \dots, Y_{ik-1}$  and  $X_i$ . This corresponds to adding a random “error” from the regression model. In addition we must account for one more source of uncertainty. Recall that the prediction of  $Y_{ik}$  is based on a set of *estimated* regression coefficients and an estimate of the residual variance,  $\sigma^2$ ; however, implicitly, the imputation process above would be treating these as fixed and known rather than sample estimates. A “proper” imputation should also account for this latter source of variability. To incorporate this additional source of variation, each imputation should be based on different (i.e., randomly perturbed) values for the regression coefficients and  $\sigma^2$ . (Specifically, these values should be random draws from what is known as their “posterior distributions”. A detailed description of “posterior distributions” is omitted because it requires some basic understanding of Bayesian statistics, a topic that is beyond the scope of this chapter). This last step, randomly drawing from the “posterior distribution” of the parameters (under a so-called non-informative prior distribution for these parameters), is probably the most

complicated part of the imputation process. Fortunately, this step has been implemented in many statistical software packages for multiple imputation (e.g., PROC MI in SAS assuming multivariate normality).

To summarize, regression methods require the use of the following two steps to produce imputed values, say  $\hat{Y}_{ik}$ , for the missing  $Y_{ik}$ :

1. New regression parameters, say  $\gamma^*$ , and the residual variance, say  $\sigma^{*2}$ , are randomly drawn from their “posterior distributions” to account for the uncertainty in estimating  $\gamma$  and  $\sigma^2$ . Specifically, the residual variance is randomly drawn as

$$\sigma^{*2} = (N_k - q) \hat{\sigma}^2 / \chi^2,$$

where  $N_k - q$  denotes the degrees of freedom for the residual variance, and  $\chi^2$  is a random draw from a chi-square distribution with  $(N_k - q)$  degrees of freedom. Then the regression parameters,  $\gamma^*$ , are randomly drawn from a multivariate normal distribution with mean equal to the estimated regression parameters,  $\hat{\gamma}$ , and with covariance matrix,

$$\widehat{\text{Cov}}(\hat{\gamma}) = \sigma^{*2} \left( \sum_{i=1}^{N_k} Z_{ik} Z'_{ik} \right)^{-1},$$

where  $\hat{\sigma}^2$  has been replaced by  $\sigma^{*2}$  and  $Z_{ik}$  denotes the design vector for the regression of  $Y_{ik}$  on  $Y_{i1}, \dots, Y_{ik-1}$ , and any subset of the components of  $X_i$ .

2. The missing values for  $Y_{ik}$ , say  $\hat{Y}_{ik}$ , can then be imputed on the basis of the following predictions:

$$Y_{ik}^* = Z'_{ik} \gamma^* + e^*,$$

where, for each missing observation on  $Y_{ik}$ ,  $e^*$  is randomly drawn from a normal distribution with mean zero and standard deviation,  $\sigma^*$ .

Regression imputation of the remaining missing values continues in a similar manner until all of the missing values have been filled in. Multiple imputations are obtained by repeating these steps  $m > 1$  times.

The description of regression imputation given above has focused on the case where the longitudinal responses are continuous. When the responses are discrete rather than continuous (e.g., repeated binary responses), a regression imputation can be based on a series of suitable generalized linear models,

$$g\{E(Y_{ik}|Y_{i1}, \dots, Y_{ik-1}, X_i)\} = Z'_{ik} \gamma,$$

where  $g(\cdot)$  is a known link function. For example, with binary responses, the missing responses can be imputed sequentially using all preceding responses in the monotone pattern (and any subset of the covariates) as predictors in a series of logistic regression models. The logistic regression model at the  $k^{th}$  occasion,

$$\text{logit}\{\Pr(Y_{ik} = 1|Y_{i1}, \dots, Y_{ik-1}, X_i)\} = Z'_{ik} \gamma,$$

can be estimated using standard statistical software for logistic regression. Based on the estimated parameters for the logistic regression model, say  $\gamma$ , new logistic regression models are obtained by randomly drawing logistic regression parameters, say  $\gamma^*$ , from their “posterior distribution.” The missing binary responses can then be imputed from Bernoulli distributions with probabilities of success,

$$\frac{\exp(Z'_{ik} \gamma^*)}{1 + \exp(Z'_{ik} \gamma^*)},$$

determined by the randomly drawn logistic regression parameters. In a similar way, when the longitudinal responses are counts, regression imputation can be based on loglinear regression models where missing responses are imputed from Poisson distributions.

# Predictive Mean Matching

The second approach to imputation is known as predictive mean matching and is closely related to regression methods. Predictive mean matching is also based on a sequence of regression models for  $Y_{ik}$ , given  $Y_{i1}, \dots, Y_{ik-1}$  and  $X_i$ , fit to the observed data. However, instead of filling in the missing values on the basis of predicted values from the regressions, predictive mean matching imputes using the observed values of the outcomes from the data at hand that are in a certain sense “closest” to the predicted values.

As with regression imputation methods, a series of regression models for  $Y_{ik}$ , given  $Y_{i1}, \dots, Y_{ik-1}$  and  $X_i$ , are fit using the observed data. When the responses are continuous, standard linear regression is typically used to obtain predicted values. The fitted linear regression produces estimates of both the regression parameters and the residual variance. New regression parameters and residual variance are then randomly drawn from their “posterior distributions” to account for the uncertainty in their estimation. Given a random draw of  $\gamma$  from its “posterior distribution”, say  $\gamma^*$ , for each missing value a predicted value

$$\hat{Y}_{ik} = Z'_{ik}\gamma^*$$

is calculated. Note that unlike in regression imputation where a random “error” is also included in the prediction, here  $\hat{Y}_{ik}$  is a prediction of the mean. Predicted values are also calculated for observations with non-missing values for  $Y_{ik}$ . From the latter, a subset of  $K$  observations whose corresponding predicted values are closest to  $\hat{Y}_{ik}$  is generated; the predictive mean matching method requires the number of closest observations,  $K$ , to be specified in advance. Thus for each missing value there is a set of  $K$  potential “donors” who have similar predicted values. The missing value is then replaced by a value drawn randomly from these  $K$  observed values. Multiple imputations are obtained by repeating these steps  $m > 1$  times.

Although predictive mean matching is closely related to regression methods, it should be somewhat less sensitive to any misspecification of the sequence of regression models. That is, predictive mean matching relies on the regression models only to determine the distance between observed and missing values. Instead of replacing the missing values using predictions from the regression models, predictive mean matching imputes using the observed values of the outcomes that are “closest” to the predicted values. Of note, because imputed values are drawn randomly from observed values contributed from other subjects, the predictive mean matching method ensures that all imputed values are plausible.

Finally, we note that there is a third approach to imputation known as “propensity score” methods. In propensity score methods, values to impute for the missing responses are randomly drawn from observations on subjects who are equally likely to drop out but do not at that occasion. Propensity score methods require a model for the probability of dropout at any occasion; ordinarily it is assumed that dropout depends only on past observed responses and any subset of the observed covariates (i.e., the dropout mechanism is MAR). However, a potential drawback of this approach is that the resulting set of imputed values is not, in general, a proper imputation of  $Y_i^M$  from  $f(Y_i^M|Y_i^O, X_i)$ . As a result this method of imputation may not preserve the inter-relationships among the variables.

## 18.2.2 Non-monotone Missing Data Patterns

When the missing data patterns are monotone it greatly simplifies the process of drawing “proper” imputations. There are many methods of imputation with monotone missingness; in the previous section we have focused on two commonly used methods. When the missing data patterns are non-monotone or intermittent, iterative and more computationally demanding methods are usually required because it is no longer straightforward to randomly sample from  $f(Y_i^M|Y_i^O, X_i)$ . In this section we focus on two commonly used algorithms for simulating random draws from  $f(Y_i^M|Y_i^O, X_i)$ . The first algorithm, known as “data augmentation” (DA), yields a random sample from  $f(Y_i^M|Y_i^O, X_i)$  after a sufficiently large number of iterations of the algorithm. The second algorithm, known as “chained equations,” is somewhat ad hoc and produces imputations that are, at best, a random sample from an approximation to  $f(Y_i^M|Y_i^O, X_i)$ . Despite the fact that the “chained equations” approach rests on shaky statistical foundations, it appears to work reasonably well in practice.

Data augmentation is best understood as an iterative two-step algorithm that alternates between random draws of the missing responses, given current values of the imputation model parameters, and random draws of the model parameters, given both the observed and current imputed values of the responses. Specifically, DA involves iterating between the following two steps:

1. randomly sampling the missing responses,  $Y_i^M$ , from  $f(Y_i^M|Y_i^O, X_i)$  given a current random draw of the imputation model parameters, and
2. randomly sampling the imputation model parameters given a current random draw of the missing data.

The resulting sequence of draws of the missing data,  $Y_i^M$ , and the imputation model parameters, form what is known as a Markov chain. Note that imputations of missing responses in step 1 are not necessarily draws from the true  $f(Y_i^M|Y_i^O, X_i)$  because they are based on imputation model parameters drawn (in step 2) from the distribution that results from treating the imputed values as if they were actual observed values. However, after a sufficiently large number of iterations of this two-step algorithm, the components of this Markov chain have the desired distribution and  $Y_i^M$  in step 1 can be considered a random sample from  $f(Y_i^M|Y_i^O, X_i)$ . Data augmentation is often referred to as a Markov chain Monte Carlo (MCMC) method because it involves repeated random sampling (also known as Monte Carlo simulation) of a Markov chain whose distribution, after a sufficiently large number of iterations, converges or stabilizes to  $f(Y_i^M|Y_i^O, X_i)$ .

When it is assumed that the responses have a multivariate normal distribution, the first step of the DA algorithm is relatively straightforward because  $f(Y_i^M|Y_i^O, X_i)$  also has a (multivariate) normal distribution. That is, given some current values of the mean vector and the covariance matrix, imputing the missing responses only requires drawing a random sample from a conditional (multivariate) normal distribution. In contrast, the random sampling in the second step is somewhat more involved. In the second step, given the observed and current imputed values for the responses (i.e., given a so-called completed data set), new model parameters must be obtained by sampling from the “posterior distribution” of the mean vector and covariance matrix. Without any prior information about the mean vector and covariance matrix, these too can be simulated from well-known distributions. Specifically, the mean vector can be randomly sampled from a multivariate normal distribution and the covariance matrix can be randomly sampled from an inverted Wishart distribution; a more detailed description of sampling from the posterior distribution of the mean vector and covariance matrix is outside the scope of this chapter. These new parameter values are then used in the first step and the process iterates, creating a Markov chain. Given a sufficiently large number of iterations, the imputed values for  $Y_i^M$  in step 1 can be considered a random sample from  $f(Y_i^M|Y_i^O, X_i)$ .

The use of MCMC methods for generating imputations requires some additional care because (1)

the Markov chain may require many iterations before the desired stationary distribution is obtained and (2) the Markov chain has an inherent dependence, in the sense that the current state of the chain has some influence on the next state in the iteration (i.e., with any Markov chain, the “current” state is predictive of the “future”). To address the first concern, a large number, say 1000 to 5000, of initial or “burn-in” iterations of the Markov chain are run, from which no imputations are made. In a sense, samples from a large number of iterations at the beginning of the algorithm are simply ignored; these burn-in iterations are executed before the first imputation is drawn from the Markov chain. The large number of “burn-in” iterations used at the start of the Markov chain increases the likelihood that the desired stationary distribution that we wish to sample from has been achieved. This also removes any dependence on the starting value selected for the chain (ordinarily determined by the computer’s random seed). To address the second concern about dependence, successive iterations of the chain are avoided as they tend to be correlated. Instead, the chain is subsampled, with imputations drawn from every  $k^{th}$  iteration of the chain (where  $k > 1$ ). By choosing a relatively large value for  $k$ , say 100 to 500, any dependence between consecutive imputations drawn from the chain is negligible.

Next we consider a second iterative method for drawing imputations when the missing data patterns are non-monotone. This alternative method does not rely on the assumption that the responses have a multivariate normal distribution. The method is referred to as “multivariate imputation by chained equations” (MICE) and requires that a sequence of separate regression models be specified for each response with missing data.<sup>1</sup> The specific type of regression model selected depends on the type of response variable; for example, a linear regression model is commonly used for a quantitative response, a logistic regression model for a binary response, and a Poisson regression model for counts. A notable feature of each of these regression models is that the response at any occasion is regressed on all other responses, both past and future responses (and any subset of the covariates). That is, MICE methods specify a sequence of regression models for  $f(Y_{ij}|Y_{i1}, \dots, Y_{i,j-1}, Y_{i,j+1}, \dots, Y_{in}, X_i)$ .

Once a sequence of regression models has been specified for the conditional mean of the response at each occasion,  $E(Y_{ij}|Y_{i1}, \dots, Y_{i,j-1}, Y_{i,j+1}, \dots, Y_{in}, X_i)$ , imputation by chained equations involves cycling through the fitting of these regression models, usually in order of the response variables that have the least amount of missing data. In each of these regression models, the response at any occasion is regressed on all other responses, both past and future (and any subset of the covariates), where any missing values among the “predictors” in the regression have been replaced by their imputed values from the previous cycle of the algorithm. After a particular regression model has been fit, imputed values for the missing responses can be generated from the fitted regression model, following the usual steps for a regression imputation to properly account for uncertainty. This sequence of imputing missing values for each response can be continued from one cycle to another, each time overwriting previously drawn values with updated imputed values. Typically this process continues iteratively for a pre-determined number of cycles (akin to “burn-in” iterations) and the set of imputations in the last cycle is used to generate a “completed” data set. Multiple imputations are obtained by repeating these steps  $m > 1$  times.

One of the appealing features of imputation by chained equations is that it avoids the problem of implicitly assuming a multivariate normal model for the responses. By specifying imputation models on a variable by variable basis, different imputation models can be used for different types of responses; e.g., a logistic regression imputation model can be used for binary responses. However, the method is somewhat ad hoc and does not have a firm theoretical basis. Specifically, the set of conditional distributions specified by each regression model,  $f(Y_{ij}|Y_{i1}, \dots, Y_{i,j-1}, Y_{i,j+1}, \dots, Y_{in}, X_i)$ , may not even be compatible with any joint distribution for the vector of responses,  $f(Y_{i1}, \dots, Y_{in}|X_i)$ . The potential consequences of this incompatibility on the quality of the imputations from MICE are unclear. So, at best, the imputations arising from this method should be regarded as random draws from an approximation to  $f(Y_i^M|Y_i^O, X_i)$ . With this caveat, there are results from limited simulation studies that suggest the method yields approximately valid inferences.

### 18.2.3 Concluding Remarks

Multiple imputation is a flexible method for handling missing data in longitudinal studies when the missing data mechanism is thought to be MAR. We note that many of the imputation models that have been reviewed in the previous sections assume a multivariate normal distribution for the response vector. Naturally, imputations based on a multivariate normal distribution for the responses will be far more appealing for the linear models for longitudinal data described in Part II of the book. For these linear models, recall that a likelihood-based analysis of the incomplete data is also valid under MAR. However, such an analysis does require that the likelihood be correctly specified. For example, maximum likelihood estimation of the regression parameters in these linear models is valid when data are MAR provided that the assumed multivariate normal distribution for the responses has been correctly specified.

When the assumptions for a likelihood-based analysis hold, multiple imputation may offer few potential benefits. In settings where the assumptions (e.g., multivariate normality) underlying both approaches are tenable, the results should be relatively similar. In such settings ML estimation based on the incomplete data may be preferred on grounds that (1) it is less computationally intensive, (2) it yields a unique set of results, and (3) it provides the most efficient estimates of the model parameters. However, conditions 2 and 3 become less important as the number of imputed values  $m$  increases. That is, as  $m$  increases (or, more formally, as  $m \rightarrow \infty$ ), multiple imputation yields almost efficient estimates when compared to ML estimation with the incomplete data; in addition, as  $m$  increases, the final estimates yielded by multiple imputation become far more stable. Of course, when  $m$  is relatively large, the computational burden of multiple imputation is far greater relative to a likelihood-based analysis.

When does multiple imputation offer distinct advantages for the analysis of longitudinal studies with missing data? There are at least three scenarios where multiple imputation may be considered advantageous. First, when there are extraneous covariates that are thought to be either predictive of the probability of missingness and/or predictive of the responses. By extraneous, we mean covariates that would not ordinarily be included in the analysis model. These extraneous covariates are sometimes referred to as auxiliary variables because there is no interest in making inferences about them (or conditional upon them). In multiple imputation, these extraneous covariates can be introduced in the imputation process in a relatively straightforward way to potentially improve the imputation of missing values. For example, in a clinical trial, missingness may be related to side effects of the treatments. In such a setting, side effects is an extraneous covariate in the sense that it would not ordinarily be included in the analysis model. That is, ordinarily there is no scientific interest in an analysis of change in the response that conditions or stratifies on whether an individual experiences side effects. However, although side effects is considered an extraneous covariate for the analysis model, it can be incorporated into the imputation model. Indeed, the inclusion of any extraneous covariate that is highly correlated with the response is likely to improve the imputations. Of course, it should be acknowledged that inclusion of these extraneous covariates in the imputation model, but not in the analysis model, implies that there is some incompatibility between the two models. In ideal circumstances the two models, the imputation model and the analysis model, should agree in terms of their representation of the relationships among the variables. However, in practice, when the analysis model is simpler than the imputation model (e.g., the analysis model implicitly assumes no dependence between the responses and extraneous covariates), this type of incompatibility should not be of great concern.

The second scenario is when there are missing covariates in addition to missing responses. Likelihood-based analysis with incomplete covariates and responses is not straightforward and has not been implemented in standard statistical software. On the other hand, multiple imputation of both missing responses and covariates is, in principle, a relatively straightforward way to handle this problem (albeit requiring some assumptions about the missing covariates).

Third, multiple imputation may be appealing in settings where a full likelihood-based analysis is not possible because there is no convenient specification of the joint multivariate distribution of the

vector of responses. For example, likelihood-based analysis of marginal models is not at all straightforward when the vector of responses is discrete rather than continuous. In such settings the generalized estimating equations (GEE) approach is routinely used as an alternative method of estimation. Standard applications of the GEE approach are valid only under MCAR. However, certain multiple imputation methods (e.g., logistic regression methods) can be fruitfully combined with standard GEE analyses of discrete longitudinal responses to make it valid under MAR. Finally, we note that multiple imputation also provides a relatively flexible and general framework for undertaking sensitivity analyses. By considering a series of alternative imputation models for the missing data, the impact of variations in the imputation model on the overall results provides an assessment of their robustness.

# 18.3 INVERSE PROBABILITY WEIGHTED METHODS

In the previous section we discussed multiple imputation methods that replace missing values with randomly drawn values from the conditional distribution of the missing data given the observed data, denoted  $f(Y_i^M|Y_i^O, X_i)$ . In this section we consider an alternative approach for handling missing data that does not require any assumptions about  $f(Y_i^M|Y_i^O, X_i)$ ; instead, an adjustment to the analysis is made by weighting the observed data in some appropriate way. In weighting methods, the under-representation of certain response profiles in the observed data is taken into account and corrected. These weighting approaches are often called propensity weighted or inverse probability weighted (IPW) methods. In general, inverse probability weighted methods are more straightforward to implement when any missingness is restricted to dropout. In addition IPW methods are more appealing in settings where a full likelihood-based analysis is not possible due to the lack of a convenient specification of the joint multivariate distribution of the vector of responses, such as in settings where the vector of responses is binary rather than continuous. As a result the following description of IPW methods focuses exclusively on the problem of handling dropout in GEE analyses of discrete longitudinal responses.

Recall that when missingness is restricted to dropout, we can replace the vector of response indicators,  $R_i = (R_{i1}, \dots, R_{in})'$ , by a scalar variable  $D_i$ , with  $D_i = 1 + \sum_{j=1}^n$  denoting the occasion at which dropout occurs. For a so-called complete-case  $Y_i = (Y_{i1}, \dots, Y_{in})'$  and  $D_i = n + 1$ , whereas for an individual with an incomplete vector of  $n_i$  responses (where  $n_i < n$ ), with observed components  $Y_i^o = (Y_{i1}, \dots, Y_{in_i})'$ ,  $D_i = n_i + 1$ .

The basic idea underlying all IPW methods is to base estimation on the observed responses but weight them to account for the inverse probability of remaining in the study. The propensities for dropout (or, conversely, for subjects remaining in the study) can be estimated as a function of the observed responses prior to dropout, and also as a function of the covariates and any additional variables (e.g., subject characteristics) that are thought likely to predict dropout. To provide some intuition for IPW methods, we first consider the simplest version of this approach based on an adjustment to the complete-case analysis. For an IPW complete-case analysis we need to estimate a single weight, say  $w_i$ , only for those subjects who complete the study ( $D_i = n + 1$ ). The weight  $w_i$  can be computed as the inverse of the product of the conditional probabilities of remaining in the study at each occasion,

$$w_i = (\pi_{i1} \times \pi_{i2} \times \dots \times \pi_{in})^{-1},$$

where  $\pi_{ij} = \Pr(D_i > j | D_i \geq j) = \Pr(R_{ij} = 1 | R_{i1} = \dots = R_{i,j-1} = 1)$ . For the first occasion, it is usually assumed that  $R_{i1} = 1$  for all individuals, and then  $\pi_{i1} = 1$ . For all subsequent occasions, the  $\pi_{ij}$ 's can be estimated from those individuals remaining in the study at the  $(j - 1)^{th}$  occasion, given their recorded history of observed responses and covariates up to the  $(j - 1)^{th}$  occasion. The weight  $w_i$  is then estimated as the inverse of the product of the estimated  $\pi_{ij}$ 's. In an IPW complete-case analysis, each contribution to the analysis from the “completers” receives a weight of  $w_i$ . The intuition behind this weighting is that each subject’s contribution to the weighted complete-case analysis is replicated  $w_i$  times, in order to count once for herself and  $(w_i - 1)$  times for those subjects with the same history of responses and covariates but who did not complete the study. For example, a subject with weight of 4 has a probability of completing the study of 0.25 (or  $\frac{1}{w_i} = 0.25$ ). As a result, in a complete-case analysis, data from this subject should count once for herself and 3 times for those subjects who do not complete the study (recall that if the probability of completing the study is  $\frac{1}{4}$ , it means that 3 subjects are expected to drop out for every one that completes the study). Given  $\hat{w}_i$ , the inverse of the estimated probability that the  $i^{th}$  subject completes the study, an IPW complete-case analysis can then be performed. For example, the standard GEE approach can be readily adapted to handle dropout

that is MAR by making an appropriate adjustment to the analysis for the propensity for dropout. Specifically, an inverse probability weighted GEE (IPW-GEE) can be used to analyze the data from the “completers” only, where each subject’s contribution to the analysis is weighted by  $\hat{w}_i$ .

An IPW complete-case analysis is valid provided that the model that produces the estimated  $w_i$  is correctly specified. The  $w_i$  are not ordinarily known when there is dropout in longitudinal studies, but must be estimated from the observed data (e.g., using a repeated sequence of logistic regressions to model the  $\pi_{ij}$ ’s). Under the assumption that data are MAR, the  $\pi_{ij}$ ’s are assumed to be a function only of the *observed* covariates and the *observed* responses prior to dropout. The  $\pi_{ij}$ ’s can also depend on any additional variables or subject characteristics that are thought likely to predict dropout. Therefore estimation of the weights is, in principle, straightforward. However, an IPW analysis restricted to the complete-cases is less than optimal in the sense that it makes very inefficient use of the available data. Next we discuss how the IPW method for handling dropout can be made more efficient by conducting an appropriately weighted available-data analysis, with weights that are also occasion-specific. Specifically, we focus on describing a general application of IPW methods to the GEE analysis of longitudinal data.

Recall from Chapter 13 that the standard GEE estimator is obtained as the solution to the following estimating equations:

$$\sum_{i=1}^N D_i^o' V_i^{o-1} (Y_i^o - \mu_i^o) = 0,$$

where  $D_i^o = \frac{\partial \mu_i^o}{\partial \beta}$ , and  $\mu_i^o = g^{-1}(X_i^o \beta)$  denotes the components of  $\mu_i$  corresponding to the observed components of the response vector  $Y_i^o = (Y_{i1}, \dots, Y_{in_i})'$ . As was discussed in Chapter 13, the solution to these estimating equations yields a consistent estimator of  $\beta$  provided the data are MCAR (or provided that missingness depends only on the covariates included in the model for the mean response). However, when dropout is MAR, the standard GEE can yield badly biased estimates of  $\beta$ . The inverse probability weighted GEE (IPW-GEE) approach was developed to circumvent this specific problem. In the IPW-GEE the dropout process is accounted for by appropriately weighting the estimating equations. Specifically, the IPW-GEE estimator is obtained as the solution to the following *weighted* estimating equations:

$$\sum_{i=1}^N D_i' V_i^{-1} W_i (Y_i - \mu_i) = 0,$$

where  $D_i$  is the  $n \times p$  derivative matrix,  $V_i$  is a  $n \times n$  working covariance matrix for  $Y_i$ , and  $W_i$  is an  $n \times n$  diagonal matrix of the occasion-specific weights,  $w_{ij}$ , for  $j = 1, \dots, n$ ,

$$W_i = \begin{pmatrix} R_{i1} \times w_{i1} & 0 & \cdots & 0 \\ 0 & R_{i2} \times w_{i2} & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & R_{in} \times w_{in} \end{pmatrix}.$$

Note that if the  $i^{th}$  individual’s response is observed at the  $j^{th}$  occasion, it receives weight of  $w_{ij}$ ; in contrast, all of the unobserved responses receive weight of zero. The weight,  $w_{ij}$ , is the inverse of the *unconditional* probability of being observed at the  $j^{th}$  occasion.

To calculate these weights, let  $\pi_{ij}$  denote the *conditional* probability of the  $i^{th}$  individual being observed (or not dropping out) at the  $j^{th}$  occasion, given that this individual was observed at the prior occasions. For the first occasion it is usually assumed that  $R_{i1} = 1$  for all individuals, and then  $\pi_{i1} = 1$ . Recall that the MAR assumption implies that

$$\begin{aligned} \pi_{ij} &= \Pr(D_i > j | D_i \geq j, X_i, Y_i) \\ &= \Pr(R_{ij} = 1 | R_{i1} = \dots = R_{i,j-1} = 1, X_i, Y_i) \\ &= \Pr(R_{ij} = 1 | R_{i1} = \dots = R_{i,j-1} = 1, X_i, Y_{i1}, \dots, Y_{i,j-1}), \end{aligned}$$

with  $\pi_{ij}$  depending only on fully-observed covariates and previously-observed responses. It is also

possible to allow the  $\pi_{ij}$  to depend on additional fully-observed variables (e.g., subject characteristics) that are thought to be predictive of dropout. For the IPW-GEE analysis the required weight  $w_{ij}$  for the  $i^{th}$  individual at the  $j^{th}$  occasion is the inverse of the *unconditional* probability of being observed at the  $j^{th}$  occasion. The *unconditional* probability of being observed at the  $j^{th}$  occasion can be expressed as the cumulative product of conditional probabilities,

$$\pi_{i1} \times \pi_{i2} \times \cdots \times \pi_{ij}.$$

The required weight is then given by the inverse of the cumulative product of conditional probabilities,

$$w_{ij} = (\pi_{i1} \times \pi_{i2} \times \cdots \times \pi_{ij})^{-1}.$$

Because we assumed that  $R_{i1} = 1$ , the first diagonal element of  $W_i$  is fixed at 1, while the remaining elements have weights given by  $w_{ij}$  for occasions when the response is observed ( $R_{ij} = 1$ ) and weights of zero for occasions when the response is not observed ( $R_{ij} = 0$ ).

Two somewhat technical, yet important, details are worth mentioning at this stage. The first relates to the magnitudes of the estimated weights. With IPW methods, extra care needs to be exercised when certain configurations of values for the “covariates” (i.e.,  $X_i$ ,  $Y_{i1}$ , ...,  $Y_{i,j-1}$ ) yield estimates of the  $\pi_{ij}$ ’s that are very small. When the estimates of  $\pi_{ij}$  are very small and close to zero, they can lead to estimated weights that are extremely large. Analyses that incorporate such extreme weights will be unduly influenced by the small subset of observations with these weights, yielding estimators of the regression parameters that are unstable and have poor precision. We recommend examining the distribution of the estimated weights for the presence of discernibly large values and checking the sensitivity of results to the inclusion of observations that receive large weights. If the results of the analysis are quite sensitive to a small number of large weights, then the IPW analysis should be reconsidered in favor of alternative methods of adjusting for missingness.

The second issue relates to the choice of working covariance matrix in the IPW-GEE. Unless the working covariance matrix,  $V_i$ , is assumed to be diagonal, the IPW-GEE requires that the covariates,  $X_i$ , are fully observed at all occasions. That is, the covariates are assumed to be known at both the occasions where the response is observed and those occasions where the response is unobserved. This will often be the case in designed studies, where the main components of  $X_i$  are treatment or exposure group indicators (i.e., time-invariant covariates), in addition to indicators of, or functions of, the intended times of measurement. However, in cases where not all components of  $X_{ij}$  are known when  $Y_{ij}$  is missing, the IPW-GEE estimator can be used if a “working independence” assumption is made. The “working independence” assumption corresponds to setting the off-diagonal elements of  $V_i$  to zero. When a “working independence” assumption is made, valid standard errors can be obtained by using the sandwich variance estimator. Furthermore we note that, when attempting to implement the IPW-GEE approach using standard statistical software for GEE, the “working independence” assumption may be required to ensure that the weights are appropriately incorporated in the analysis. For example, occasion-specific weights can be specified by using the WEIGHT statement in PROC GENMOD in SAS or the pweight option for the glm command in Stata; however, these weights are appropriately incorporated in the IPW-GEE analysis only under a “working independence” assumption.

Thus an important property of the IPW-GEE is that the choice of working covariance matrix only has an impact on the efficiency of estimation. The IPW-GEE does not require correct specification of the working covariance matrix to consistently estimate the components of  $\beta$  and their standard errors. It does, however, require correct specification of the model for the dropout process, that is, for valid estimation of  $\beta$ , it requires that the model that produces the estimated weights be correctly specified.

Next we consider estimation of the weights. Recall that the MAR assumption implies that  $\pi_{ij}$  depends only on *observed* covariates and *observed* past responses. Therefore we can estimate  $\pi_{ij}$  (for  $j > 1$ ) by, for example, constructing a logistic regression model for  $\pi_{ij}$ ,

$$\begin{aligned}
\text{logit}(\pi_{ij}) &= \text{logit}\{\Pr(R_{ij} = 1 | R_{i1} = \dots = R_{i,j-1} = 1, Z_{ij})\} \\
&= \theta_1 Z_{ij1} + \theta_2 Z_{ij2} + \dots + \theta_q Z_{ijq} \\
&= Z'_{ij} \theta,
\end{aligned}$$

where  $Z_{ij}$  is a  $q \times 1$  design vector that incorporates certain components of  $X_{ij}$  and the past responses ( $Y_{i1}, \dots, Y_{i,j-1}$ ), and perhaps indicator or dummy variables for occasions. Note that  $Z_{ij}$  can also incorporate additional covariates that may be predictive of dropout but are not of subject-matter interest in the marginal model for the mean response; this is an especially appealing feature of IPW methods. Estimates of the logistic regression parameters,  $\theta$ , can be obtained using standard statistical software for fitting logistic regression. That is, we can create a “stacked” data set in which each individual contributes a sequence of binary “outcomes” to the analysis, where each binary “outcome,”  $R_{ij}$ , is an indicator of whether the response was observed at a given occasion, from the second occasion (because it is assumed that  $R_{i1} = 1$ ) until either the occasion when dropout occurs or the last intended measurement occasion. Thus, individuals who do not dropout contribute a sequence of  $n - 1$  binary responses ( $R_{i2}, \dots, R_{in}$ , where  $R_{i2} = \dots = R_{in} = 1$ ) to the logistic regression analysis, whereas an individual who drops-out at the  $k^{th}$  occasion contributes  $k - 1$  binary responses ( $R_{i2}, \dots, R_{ik}$ , where  $R_{i2} = \dots = R_{i,k-1} = 1, R_{ik} = 0$ ). The covariates,  $Z_{ij}$ , in the logistic regression model can include certain components of  $X_{ij}$  and the previous observed responses ( $Y_{i1}, \dots, Y_{i,j-1}$ ), and perhaps indicator variables for occasions. The logistic regression analysis can also include additional covariates that are thought to be predictive of dropout but are not incorporated in  $X_{ij}$ . A standard logistic regression analysis of the “stacked” data set (often referred to as a “pooled” logistic regression) provides estimates of  $\theta$ ; the  $\pi_{ij}$ ’s and the required weight  $w_{ij}$  can then be estimated as

$$\hat{w}_{ij} = (\hat{\pi}_{i1} \times \hat{\pi}_{i2} \times \dots \times \hat{\pi}_{ij})^{-1}.$$

Once the set of weights have been determined, they can be incorporated in a relatively straightforward way into standard statistical software for longitudinal analyses, for example, using the WEIGHT statement within PROC GENMOD in SAS or the pweight option for the `glm` command in Stata. The resulting IPW-GEE estimates are unbiased for the regression parameters under a MAR dropout process, provided the model for estimating the dropout probabilities has been correctly specified. When the weights are incorporated in this manner, the standard errors are estimated by implicitly assuming the weights are fixed and known. In principle, it is possible to adjust the standard errors for the estimation of these weights. In practice, however, it is not straightforward to implement such an adjustment to the standard errors within existing statistical software for GEE. The formula for making an adjustment to the standard errors is somewhat involved; for completeness, it is outlined in a separate section at the end of this chapter that can be omitted at first reading without loss of continuity. Counter-intuitively, failure to account for the estimation of the weights will, in general, result in standard errors that are too large (i.e., estimation of the weights from the data at hand leads to improvements in precision). Therefore the unadjusted standard errors (e.g., based on the conventional sandwich variance estimator) provide valid inferences and can be considered to be slightly conservative. By “conservative”, we mean that nominal  $p$ -values may be slightly larger and confidence intervals may be slightly wider than they should be. Until such time as the IPW-GEE with adjusted standard errors is implemented in widely available software, we recommend basing inferences on the unadjusted standard errors.

## 18.4 CASE STUDIES

In these case studies we illustrate some of the methods described earlier for handling missing data using two examples. The first example is based on data from the Treatment of Lead-Exposed Children (TLC) Trial (see Section 5.4). Recall that in this study the response variable, blood lead levels, is continuous and measured repeatedly at four occasions. Although the data on blood lead levels for the 100 children from the succimer and placebo groups are complete, for pedagogical purposes we created an incomplete data set with a non-monotone pattern of missing values generated under a MAR mechanism. The second example is from a longitudinal clinical trial of contracepting women (Machin et al., 1988) discussed in Sections 14.7 and 15.5. In this trial the response variable is binary, indicating whether a women experienced amenorrhea in four successive injection intervals. This trial had substantial dropout.

# Treatment of Lead-Exposed Children (TLC) Trial

In our first illustration, we focus on methods for handling missing data on a continuous response. Recall that the TLC trial was a placebo-controlled, randomized trial of an orally administered chelating agent, succimer, in children with confirmed blood lead levels of 20 to 44  $\mu\text{g}/\text{dL}$ . The following analyses are based on data on blood lead levels at baseline (or week 0), week 1, week 4, and week 6 for 100 children from the succimer and placebo groups of the TLC trial.

In Chapter 5 the results of an analysis of response profiles of complete data on these 100 children were presented (see Section 5.4). The REML estimates of the regression parameters, and their standard errors, are reproduced in [Table 18.1](#). For ease of interpretation, baseline (week 0) is chosen as the reference level for time and the placebo group is chosen as the reference level for treatment. In the TLC trial the question of main scientific interest concerns the comparison of the two treatment groups in terms of their mean changes from baseline. This question translates directly into a test of the three single-degree-of-freedom contrasts for the group  $\times$  time interaction. The results in [Table 18.1](#) indicate that children treated with succimer have a discernibly greater decrease in mean blood lead levels from baseline at all occasions when compared to the children treated with placebo. For example, when compared to the placebo group, the succimer group has an additional 3.152  $\mu\text{g}/\text{dL}$  (with SE = 1.257) decrease in mean blood lead levels from baseline to week 6. There are even larger differences between the two treatment groups earlier in the trial.

**Table 18.1** Estimated regression coefficients and standard errors based on an analysis of response profiles of the blood lead level data at baseline, week 1, week 4, and week 6 for the children from the TLC trial.

Variable	Group	Week	Estimate	SE	Z
Intercept			26.272	0.710	36.99
Group	S		0.268	1.005	0.27
Week		1	-1.612	0.792	-2.04
Week		4	-2.202	0.815	-2.70
Week		6	-2.626	0.889	-2.96
Group $\times$ Week	S	1	-11.406	1.120	-10.18
Group $\times$ Week	S	4	-8.824	1.153	-7.66
Group $\times$ Week	S	6	-3.152	1.257	-2.51

To create an incomplete data set, we applied the following missing at random (MAR) mechanism to the responses at weeks 1, 4, and 6,

$$\text{logit}\{\Pr(R_{ij} = 1|Y_{i1}, \text{Group}_i)\} = \theta_1 + \theta_2 Y_{i1} \times \text{Group}_i \times (j - 1), \quad j = 2, 3, 4.$$

When  $\theta_2 = 0$ , this mechanism is MCAR with a constant probability of missingness,  $\frac{1}{1+\exp(\theta_1)}$ . However, when  $\theta_2 \neq 0$ , missingness is allowed to depend on the baseline response ( $Y_{i1}$ ) for children randomized to the succimer group ( $\text{Group}_i = 1$ ), with the strength of that dependence increasing over time. Note that this missing data mechanism yields a non-monotone pattern of missingness. To generate an incomplete data set, we fixed  $\theta_1 = 2.5$  and  $\theta_2 = -0.03$ ; a negative value for  $\theta_2$  implies that children in the succimer group with higher blood lead levels at baseline have a greater probability of being missing at subsequent occasions. This yielded the non-monotone patterns of missingness displayed in [Table 18.2](#). In the placebo group, 6% of children have missing responses at week 1, 8% at week 4, and 4% at week 6; in the succimer group, the corresponding rates are 18% at week 1, 30% at week 4, and 48% at week 6. Because missingness is related to higher baseline blood lead levels in the succimer group, we might expect the mean blood lead levels in the succimer group to be lower, when compared to the means in the complete data set, at subsequent occasions. The discrepancy between the means in the complete and incomplete data sets should be most pronounced at week 6 when the rate of missingness in the succimer group is relatively high. This trend is

apparent in the sample means reported in [Table 18.3](#), where the mean in the succimer group at week 6, based on the incomplete data, is approximately 1.5 units lower than the corresponding mean in the complete data set.

**Table 18.2** Missing data patterns generated in the succimer and placebo groups from the TLC trial.

Group	Week 0	Week 1	Week 4	Week 6	Frequency	Percent
Succimer	O	O	O	O	18	36%
	O	O	O	M	12	24%
	O	O	M	O	5	10%
	O	O	M	M	6	12%
	O	M	O	O	2	4%
	O	M	O	M	3	6%
	O	M	M	O	1	2%
	O	M	M	M	3	6%
Placebo	O	O	O	O	42	84%
	O	O	O	M	2	4%
	O	O	M	O	3	6%
	O	M	O	O	2	4%
	O	M	M	O	1	2%

Note: O denotes observed response, M denotes missing response.

**Table 18.3** Mean blood lead levels at baseline, week 1, week 4, and week 6 for the incomplete data on children from the TLC trial. Mean blood lead levels for the complete data are reported in parentheses.

Group	Baseline	Week 1	Week 4	Week 6
Succimer	26.5	13.4	15.2	19.3
	(26.5)	(13.5)	(15.5)	(20.8)
Placebo	26.3	24.5	24.1	23.7
	(26.3)	(24.7)	(24.1)	(23.6)

Next we compare and contrast a number of alternative methods for analyzing the incomplete data set on blood lead levels. First we consider a likelihood-based analysis of response profiles (see Sections 5.1–5.3) of the incomplete data. Recall that maximum likelihood (ML) estimation of the regression coefficients is valid when data are MAR provided that the assumed multivariate normal distribution for the responses has been correctly specified. This requires correct specification of not only the model for the mean response but also the model for the covariance among the responses. Because the analysis of response profiles gives unrestricted estimates of the means in each group, and also assumes an unstructured covariance among the responses, a likelihood-based analysis of the incomplete data should yield unbiased estimates of the regression coefficients. The results of an analysis of the incomplete data, summarized in terms of the three single-degree-of-freedom contrasts for the group  $\times$  time interaction, are presented in [Table 18.4](#); for ease of comparison the corresponding results for the complete data are also reproduced in [Table 18.4](#). For simplicity we focus on the treatment group effect on changes in the mean blood lead levels from baseline to week 6. With almost 50% of the responses missing at week 6, this is the effect most likely to be sensitive to missingness. The analysis of response profiles of the incomplete data yields an estimate of  $-2.937$  (with SE = 1.532) for the treatment group effect on changes in the mean blood lead levels from baseline to week 6. This is quite similar to the corresponding estimate from the complete data, albeit with a standard error that is discernibly larger. The larger standard error is to be expected and reflects the loss of information due to the relatively large fraction of missing data at week 6. The estimates of the remaining two contrasts for the group  $\times$  time interaction are also quite similar to those obtained from the complete data. In general, we would expect to obtain similar estimates of effects under the assumption that missingness is MAR and that the models for both the mean and the covariance have been correctly specified. In this case the former assumption is known to be valid by

definition of the missing data mechanism that created the incomplete data set. The latter assumption seems tenable given that the analysis of response profiles makes few assumptions about the structure of the mean response and the covariance among the responses.

**Table 18.4** Estimated regression coefficients (and standard errors) based on an analysis of response profiles of (i) the complete data, denoted ML(C), (ii) the incomplete data, under the assumption of an unstructured covariance, denoted ML(UN), (iii) the incomplete data, under the (incorrect) assumption of independent responses and constant variance, denoted ML(IND), and (iv) multiple imputed data sets, denoted ML(MI).

Variable	Group	Week	Complete	Incomplete		
			ML(C)	ML(UN)	ML(IND)	ML(MI)
Group × Week	S	1	−11.406 (1.120)	−11.276 (1.213)	−11.372 (1.325)	−11.279 (1.216)
Group × Week	S	4	−8.824 (1.153)	−8.985 (1.189)	−9.120 (1.342)	−8.986 (1.210)
Group × Week	S	6	−3.152 (1.257)	−2.937 (1.532)	−4.669 (2.021)	−3.030 (1.542)

It is instructive to examine the sensitivity of ML estimation with incomplete data to misspecification of the likelihood. Specifically, we consider an analysis of response profiles under the naive assumption of constant variance and independence among the repeated responses, that is, under misspecification of the covariance. This analysis yields an estimate of −4.669 (with SE = 2.021) for the treatment group effect on changes in the mean blood lead levels from baseline to week 6; note, to correct for misspecification of the covariance, the standard errors reported in [Table 18.4](#) are based on the empirical (or “sandwich”) variance estimator. The estimate of −4.669 is approximately 60% larger (in absolute value) than the estimate obtained assuming an unstructured covariance; it is also approximately 50% larger than the estimate obtained from the complete data. This reflects the bias that can be introduced in a likelihood-based analysis of incomplete data when the model for the covariance is not correctly specified, even when missingness is truly MAR. As was noted earlier (see Section 18.1), when there are missing data, it is important to correctly specify the models for both the mean and the covariance to ensure that the regression parameters for the mean model are estimated without bias due to missingness.

A final approach to the analysis of the incomplete data is to use multiple imputation to replace any missing values by a set of  $m$  plausible values. Because the missingness patterns are non-monotone and the response is continuous, we use Markov chain Monte Carlo (MCMC) methods, based on a multivariate normality assumption, to impute missing values. Specifically, missing values are replaced by a set of  $m = 50$  imputed values sampled from a chain that first uses 5000 “burn-in” iterations to achieve the desired stationary distribution. That is, 5000 burn-in iterations are executed before the first imputation from the chain is obtained. In addition, to remove any dependence between consecutive imputations, samples from the stationary distribution are drawn 50 times at every subsequent 500 iterations of the chain. Results of analyses of response profiles of the 50 completed data sets are then appropriately combined to yield the estimates and standard errors reported in the last column of [Table 18.4](#).

The analysis based on multiple imputation yields an estimate of −3.030 (with SE = 1.542) for the treatment group effect on changes in the mean blood lead levels from baseline to week 6. This estimate and standard error are very similar to those yielded by the likelihood-based analysis of response profiles of the incomplete data. In general, we might expect the analysis based on multiple imputation to yield similar results to those from a likelihood-based analysis of the incomplete data because both sets of analyses are based on the assumption of multivariate normality.

# Clinical Trial of Contracepting Women

Next we illustrate some of the methods described earlier for handling dropouts when the response is categorical rather than continuous. The methods are applied to data on a binary response from the longitudinal clinical trial of contracepting women (Machin et al., 1988) discussed in Sections 14.7 and 15.5. Recall that the goal of this trial was to compare the two treatments (100 mg or 150 mg of DMPA) in terms of how the rates of amenorrhea change over time with continued use of the contraceptive method. That is, the main interest is in an analysis that compares the rates of amenorrhea over time if those women who dropped out had remained on their assigned treatment. This is sometimes called an *explanatory* analysis (Schwartz and Lellouch, 1967). An “explanatory analysis,” often referred to as an “as treated” analysis, focuses on what is thought to be the true underlying biological effects of the different treatments.

In this clinical trial a total of 1151 women completed menstrual diaries, and the diary data were used to generate a binary sequence for each woman, indicating whether or not she had experienced amenorrhea in four successive intervals. A feature of this trial is that there was substantial dropout. When the dropout rates are broken down by dosage group, the rates were marginally higher in the 150 mg dose group. Among women randomized to 100 mg (150 mg) of DMPA, 37% (39%) dropped out before the completion of the trial, 17% (17%) dropped out after receiving only one injection of DMPA, 12% (15%) dropped out after receiving only two injections, and 8% (6%) dropped out after receiving three injections. For women who dropped out before the end of the 3-month interval between injections, a determination of whether or not they experienced amenorrhea was made, on a proportionate basis, using their existing menstrual diary data for that interval.

Letting  $Y_{ij} = 1$  if the  $i^{th}$  woman experienced amenorrhea in the  $j^{th}$  injection interval, we consider the following logistic regression model for the marginal probability of amenorrhea:

$$\text{logit}(\mu_{ij}) = \beta_1 + \beta_2 t_{ij} + \beta_3 t_{ij}^2 + \beta_4 \text{Dose}_i + \beta_5(t_{ij} \times \text{Dose}_i) + \beta_6(t_{ij}^2 \times \text{Dose}_i),$$

where  $\mu_{ij} = \Pr(Y_{ij} = 1)$ ,  $t = 0, 1, 2, 3$  for the four consecutive 3-month injection intervals,  $\text{Dose} = 1$  if randomized to 150 mg of DMPA, and  $\text{Dose} = 0$  otherwise. First we consider complete-case and available-data analyses of the data. If dropout is completely at random, then valid estimates of the marginal regression parameters can be obtained using a standard generalized estimating equations (GEE) approach. To account for the within-subject association among the repeated measures, we fit six separate pairwise log odds ratios. We note that the empirical and model-based standard errors are very similar in all of the analyses. For illustrative purposes only, we also consider a GEE analysis using LVCF imputation of the missing data. The differences in results are more easily discerned by considering the dose-specific estimated rates of amenorrhea in each of the injection intervals for the quadratic trend model given above (see [Table 18.5](#)).

**Table 18.5** Estimated marginal rates of amenorrhea for quadratic trend model using GEE under three different methods for handling dropouts: complete-case (CC), available-data (AD), and last value carried forward (LVCF).

Method	Time	100 mg	150 mg	Difference	SE	Z
CC	3 months	0.176	0.155	-0.021	0.027	-0.79
	6 months	0.258	0.317	0.059	0.028	2.07
	9 months	0.369	0.463	0.094	0.033	2.83
	12 months	0.502	0.540	0.038	0.037	1.03
AD	3 months	0.184	0.201	0.017	0.023	0.73
	6 months	0.274	0.363	0.089	0.025	3.55
	9 months	0.388	0.499	0.111	0.030	3.68
	12 months	0.517	0.572	0.055	0.036	1.52
LVCF	3 months	0.184	0.201	0.017	0.023	0.75
	6 months	0.263	0.344	0.081	0.024	3.43
	9 months	0.350	0.453	0.103	0.027	3.78
	12 months	0.437	0.498	0.061	0.029	2.10

Overall, the results of the complete-case and available-data GEE analyses suggest that the rates of amenorrhea in the second and third injection intervals are significantly higher for women who received the higher dose of DMPA, although these differences tend to decline by the end of the study. For example, during the third injection interval (6–9 months post-randomization) the predicted rates of amenorrhea from the available-data analysis are 0.499 in the 150 mg dose group and 0.388 in the 100 mg dose group. However, by the final follow-up visit there is no longer a discernible treatment difference, with predicted rates of amenorrhea of 0.572 in the 150 mg dose group and 0.517 in the 100 mg dose group.

The GEE analysis based on LVCF imputation produces discernibly lower estimated rates of amenorrhea during the third and fourth intervals, when compared to the available-data analysis, although the estimates of the treatment comparisons are not too dissimilar; however, the latter cannot be expected in general. Because LVCF uses a single imputation and does not reflect any uncertainty in the imputation, the standard errors for the estimated treatment comparisons are too small. Consequently, in contrast to the other methods, the analysis based on LVCF suggests that there are treatment differences in the estimated rates of amenorrhea at the end of the trial.

Note that if dropout is not completely at random, the complete-case and available-data GEE analyses of these data can yield biased estimates of the effects of treatment. Next we consider handling dropout using inverse probability weighted and multiple imputation methods. Recall that inverse probability weighted methods require a model for the probability of dropout. We considered the following logistic regression model that assumes the log odds of remaining in the study (or, conversely, of dropout) depends on the previous observed response. Specifically, the model for being observed at the  $j^{th}$  occasion is given by

$$\text{logit}(\pi_{ij}) = \theta_1 + \theta_2 I(t = 2) + \theta_3 I(t = 3) + \theta_4 \text{Dose}_i + \theta_5 Y_{ij-1} + \theta_6 (\text{Dose}_i \times Y_{ij-1}),$$

where  $\pi_{ij} = \Pr(R_{ij} = 1 | R_{i1} = \dots = R_{ij-1} = 1, Y_{ij-1}, \text{Dose}_i)$ . In this model the log odds of remaining in the study is allowed to vary over measurement occasions, to depend on dose group, and to depend on the previous observed response; the latter dependence is also allowed to vary by dose group.

By creating a “stacked” data set, estimates of the logistic regression parameters,  $(\theta_1, \dots, \theta_6)$ , can be obtained using standard statistical software for fitting logistic regression. That is, we can create a “stacked” data set in which each individual contributes a sequence of binary “outcomes” to the logistic regression analysis. In this analysis each binary “outcome,”  $R_{ij}$ , indicates whether the response was observed at a given occasion, from the second occasion (because  $R_{i1} = 1$  for all individuals) until either the occasion when dropout occurred or the last intended measurement occasion. Thus study “completers” contribute a sequence of three binary responses ( $R_{i2} = R_{i3} = R_{i4} =$

1) to the logistic regression analysis. Individuals who dropped out at the fourth occasion also contribute three binary responses ( $R_{i2} = R_{i3} = 1, R_{i4} = 0$ ) to the analysis. In contrast, individuals who dropped out at the third occasion contribute only two binary responses ( $R_{i2} = 1, R_{i3} = 0$ ), while individuals who dropped out at the second occasion contribute a single binary response ( $R_{i2} = 0$ ).

The results of fitting the logistic regression model to this stacked data set are presented in [Table 18.6](#). There is strong evidence that the probability of dropout is related to the previous response, although this dependence does not vary significantly between the dose groups. Specifically, for individuals in the 100 mg dose group, the conditional odds of dropout is approximately 60% higher ( $\exp(0.451) = 1.57$ ) if they experienced amenorrhea at the previous occasion. Similarly, for those in the 150 mg dose group, the conditional odds of dropout is approximately two times higher ( $\exp(0.451 + 0.238) = 1.99$ ) if they experienced amenorrhea at the previous occasion. The estimated logistic regression coefficients can be used to obtain the estimated conditional probability of remaining in the study at each occasion for each individual,  $\hat{\pi}_{ij}$ . The required weight  $w_{ij}$  for the  $i^{th}$  individual at the  $j^{th}$  occasion is then estimated by

$$\hat{w}_{ij} = (\hat{\pi}_{i1} \times \hat{\pi}_{i2} \times \cdots \times \hat{\pi}_{ij})^{-1}.$$

**Table 18.6** Estimated regression coefficients and standard errors from logistic regression model for the probability of remaining in the study.

Variable	Estimate	SE	Z
Intercept	1.668	0.104	15.98
I( $t=2$ )	0.137	0.119	1.15
I( $t=3$ )	0.729	0.144	5.06
Dose	0.068	0.131	0.52
$Y_{ij-1}$	-0.451	0.162	-2.79
Dose $\times Y_{ij-1}$	-0.238	0.220	-1.08

Because the first response was fully observed, with  $R_{i1} = 1$  for all individuals,  $\pi_{i1} = 1$  by definition. Thus, at the first occasion the weight is fixed at 1 for all individuals, while at all subsequent occasions the weights are given by  $\hat{w}_{ij}$  for occasions when the response is observed ( $R_{ij} = 1$ ) and weights of zero for occasions when the response is not observed ( $R_{ij} = 0$ ). Prior to conducting an IPW-GEE analysis, we examined the distribution of the estimated weights for the presence of discernibly large weights. The estimated weights ranged from 1.0 to 2.1, so there was no concern that a small subset of the observations might have undue influence on the analysis.

To conduct IPW-GEE estimation of the logistic regression model for the marginal probability of amenorrhea,

$$\text{logit}(\mu_{ij}) = \beta_1 + \beta_2 t_{ij} + \beta_3 t_{ij}^2 + \beta_4 \text{Dose}_i + \beta_5 (t_{ij} \times \text{Dose}_i) + \beta_6 (t_{ij}^2 \times \text{Dose}_i),$$

these occasion-specific weights need to be incorporated in the analysis. This can be achieved, for example, using the WEIGHT statement within PROC GENMOD in SAS or the pweight option for the glm command in Stata. However, to ensure that the weights are appropriately incorporated in the IPW-GEE analysis, it is necessary to make a “working independence” assumption for the within-subject association among the responses. Because a “working independence” assumption is made, standard errors are based on the “sandwich” variance estimator. The results of the IPW-GEE analysis are presented in [Table 18.7](#). As in [Table 18.5](#), the results are expressed in terms of the dose-specific estimated rates of amenorrhea in each of the injection intervals from the fitted quadratic trend model. The results from the IPW-GEE analysis are very similar to those obtained from the available-data analysis. Both the point estimates of the rates of amenorrhea and their standard errors are similar. The treatment group difference in the rates of amenorrhea at the fourth injection interval yielded by the IPW-GEE analysis is marginally lower than the corresponding difference obtained from the available-data analysis (4.9% versus 5.5%, respectively). Under the assumption that

dropout is at random, and the model for dropout has been correctly specified, the results of the IPW-GEE analysis suggest that the rates of amenorrhea in the second and third injection intervals are significantly higher for those women who received the higher dose of DMPA, although these differences tend to decline by the end of the study.

**Table 18.7** Estimated marginal rates of amenorrhea for quadratic trend model using GEE under two different methods for handling dropouts: inverse probability weighted (IPW) GEE, and multiple imputation (MI) using logistic regression.

Method	Time	100 mg	150 mg	Difference	SE	Z
IPW	3 months	0.183	0.200	0.017	0.023	0.73
	6 months	0.276	0.361	0.084	0.026	3.29
	9 months	0.393	0.496	0.104	0.031	3.34
	12 months	0.521	0.570	0.049	0.037	1.33
MI	3 months	0.183	0.199	0.016	0.026	0.60
	6 months	0.276	0.365	0.089	0.026	3.43
	9 months	0.391	0.504	0.113	0.030	3.82
	12 months	0.518	0.579	0.061	0.033	1.85

Finally, we consider an analysis of the amenorrhea data that handles dropout using multiple imputation. Because the response is binary, imputation methods that are either explicitly or implicitly based on a multivariate normal distribution assumption for the response vector are unappealing. Instead, we can use a logistic regression imputation. Specifically, we consider a sequence of logistic regression models at the second through fourth occasion that assume the log odds of amenorrhea at each occasion to depend on all past observed responses, dose group, and their interactions. For example, the model at the third occasion is

$$\begin{aligned} \text{logit} \{ \Pr(Y_{i3} = 1 | Y_{i1}, Y_{i2}, \text{Dose}_i) \} &= \theta_1 + \theta_2 \text{Dose}_i + \theta_3 Y_{i1} + \theta_4 Y_{i2} \\ &\quad + \theta_5 \text{Dose}_i \times Y_{i1} + \theta_6 \text{Dose}_i \times Y_{i2}. \end{aligned}$$

Three separate logistic regression models are fit to the data at the second through fourth occasion. Based on the estimated model parameters at each occasion, new logistic regression models are obtained by randomly drawing logistic regression parameters from their “posterior distribution.” Finally, the missing binary responses at that occasion are then imputed from Bernoulli distributions with probabilities of amenorrhea determined by the randomly drawn logistic regression parameters. Logistic regression imputation of the remaining missing values continues in a similar manner until a completed data set has been created. To create 25 imputations, these steps are repeated 25 times, and results from the 25 analyses of the filled-in data sets are appropriately combined.

The results from the analysis based on multiple imputation (see bottom of [Table 18.7](#)) are similar to those obtained from the available-data analysis (see [Table 18.5](#)); they are also similar to those obtained from the IPW-GEE analysis. The treatment group difference in the rates of amenorrhea at the fourth injection interval yielded by the multiple imputation analysis is marginally higher than the corresponding difference obtained from the IPW-GEE analysis (6.1% versus 5.5%, respectively). This confluence of the estimates from the available-data, IPW-GEE, and multiple imputation analyses provides some degree of reassurance that the main conclusions of the treatment group comparisons are not very sensitive to the methods used to handle dropout. The results from the multiple imputation analysis confirm the earlier findings that the rates of amenorrhea in the second and third injection intervals are significantly higher for those women who received the higher dose of DMPA, although these differences tend to decline by the end of the study.

# 18.5 “SANDWICH” VARIANCE ESTIMATOR ADJUSTING FOR ESTIMATION OF WEIGHTS\*

In earlier sections we mentioned that the weights used in IPW-GEE are not known but must be estimated from the data at hand. It is possible to adjust the conventional standard errors for  $\hat{\beta}$ , based on the “sandwich” variance estimator, to account for the estimation of the weights used by IPW-GEE. Indeed, such an adjustment will, in general, result in smaller standard errors. In this section we briefly outline how to construct the “sandwich” variance estimator for IPW-GEE that adjusts for estimation of the weights.<sup>†</sup>

Recall that the IPW-GEE estimator of  $\beta$  is obtained as the solution to the following weighted estimating equations:

$$S(\beta) = \sum_{i=1}^N S_i(\beta) = \sum_{i=1}^N D_i' V_i^{-1} W_i (Y_i - \mu_i) = 0.$$

If the estimation of the weights in  $W_i$  is completely ignored, and  $W_i$  is assumed fixed and known, then the usual formula for the “sandwich” variance estimator is

$$\text{Cov}(\hat{\beta}) = \left( \sum_{i=1}^N D_i' V_i^{-1} W_i D_i \right)^{-1} \left( \sum_{i=1}^N S_i(\beta) S_i'(\beta) \right) \left( \sum_{i=1}^N D_i' V_i^{-1} W_i D_i \right)^{-1}.$$

In practice, a logistic regression analysis (or any other suitable model) provides estimates of  $\theta$  and the  $\pi_{ij}$ 's; thus the weights actually used in IPW-GEE are *estimated* rather than known. Specifically, the estimates of  $\theta$  (and the  $\pi_{ij}$ 's), and hence the weights  $w_{ij}$ , are obtained as the solution to the following estimating equations:

$$S(\theta) = \sum_{i=1}^N S_i(\theta) = \sum_{i=1}^N \sum_{j=2}^n R_{i,j-1} Z_{ij} (R_{ij} - \pi_{ij}) = 0.$$

These are the equations for the logistic regression analysis of the “stacked” data set. The solution to these estimating equations provides estimates of  $\theta$  and the  $\pi_{ij}$ 's; the required weights  $w_{ij}$  are then estimated as

$$\hat{w}_{ij} = (\hat{\pi}_{i1} \times \hat{\pi}_{i2} \times \cdots \times \hat{\pi}_{ij})^{-1}.$$

The IPW-GEE estimate of  $\beta$  is then obtained using *estimated* weights, denoted by  $W_i(\hat{\theta})$ . In large samples the resulting IPW-GEE estimator has a multivariate normal distribution, with mean  $\beta$ , and estimator of its covariance given by the *adjusted* “sandwich” variance estimator (where the adjustment is for estimation of the weights),

$$\text{Cov}(\hat{\beta}) = \left( \sum_{i=1}^N D_i' V_i^{-1} W_i(\hat{\theta}) D_i \right)^{-1} \left( \sum_{i=1}^N S_i^*(\beta, \theta) S_i^{**}(\beta, \theta) \right) \left( \sum_{i=1}^N D_i' V_i^{-1} W_i(\hat{\theta}) D_i \right)^{-1},$$

where

$$S_i^*(\beta, \theta) = S_i(\beta) - \left( \sum_{i=1}^N S_i(\beta) S_i'(\theta) \right) \left( \sum_{i=1}^N S_i(\theta) S_i'(\theta) \right)^{-1} S_i(\theta).$$

Note that  $S_i^*(\beta, \theta)$  is the residual from a multivariate linear regression of  $S_i(\beta)$  on  $S_i(\theta)$ . The replacement of  $S_i(\beta)$  by  $S_i^*(\beta, \theta)$  in the “meat” (or center-piece) of the “sandwich” variance estimator yields an appropriate adjustment to the standard errors of  $\hat{\beta}$  for estimation of the weights in the IPW-GEE method. Moreover, because  $S_i^*(\beta, \theta)$  is a residual, by definition, the “residual sums of squares” in the “meat” of this “sandwich” variance estimator must be smaller than (or equal to) the corresponding “sums of squares” based on  $S_i(\beta)$  in the conventional “sandwich” variance estimator. This “sums of squares” argument helps explain why the *adjusted* “sandwich” variance estimator yields smaller standard errors for  $\hat{\beta}$ .

## 18.6 COMPUTING: MULTIPLE IMPUTATION USING PROC MI IN SAS

As described in earlier sections of this chapter, making inferences from a statistical analysis based on multiple imputation is a three-step process. First, the missing data are filled-in  $m$  times to create  $m$  completed data sets. Second, the  $m$  completed data sets are analyzed using standard statistical methods (e.g., fitting linear models using PROC MIXED in SAS). Finally, the results from the  $m$  analyses of the completed data sets are appropriately combined. In this section we focus on the first step and briefly discuss how to create  $m$  completed data sets using PROC MI in SAS. There is a complementary procedure in SAS, PROC MIANALYZE, for combining the results from the  $m$  analyses of the completed data sets.

The MI procedure in SAS is a general procedure for multiple imputation that offers a number of alternative methods for imputing missing data, especially when the patterns of missing data are monotone. No attempt is made here to give a comprehensive review of the main features of PROC MI. Instead, we present illustrative source code for generating  $m$  completed data sets using a number of alternative methods for imputation.

Before discussing the command syntax for PROC MI, we note that when imputing longitudinal data with missing responses, it is important to structure the data set appropriately. To capitalize on the correlation among the repeated measurements of the responses, the procedure requires that each repeated measurement in a longitudinal data set be a separate variable rather than a separate “record.” That is, for the purposes of imputing missing values, PROC MI should be applied to a data set that is structured in a “wide” rather than a “long” format, with a single “record” for each individual. In a “wide” format, the imputation model for a missing response at any particular occasion can include as predictors the responses at any of the remaining occasions, thereby capitalizing on the positive correlation among the repeated measurements. After  $m$  completed data sets have been generated, each of these “wide” format data sets can then be transformed to a “long” format data set prior to analysis using standard statistical methods (e.g., PROC MIXED or PROC GENMOD in SAS). Finally, we note that when using imputation methods that rely on a multivariate normal assumption, imputations can often be improved by transforming the response variable. For example, when the longitudinal responses are quantitative but have distributions that are not symmetric, transformation of the response (e.g., log transformation of variables with positive skewed distributions) should improve the imputation. After the values have been imputed on the transformed scale, these can be transformed back to the original scale (e.g., if a response variable has been log transformed prior to imputation, the imputed values can subsequently be exponentiated). Of course, the multivariate normal assumption has no consequence for those variables that have no missing data but are included in the imputation process (e.g., treatment or exposure group indicators).

To use multiple imputation to generate 25 completed data sets in the setting where there is intermittent (non-monotone) missingness and the vector of response is assumed to have a multivariate normal distribution, we can use the illustrative SAS commands given in [Table 18.8](#). By default, PROC MI in SAS uses MCMC methods for generating imputations when missingness is non-monotone. In contrast, when the missing data patterns are monotone there are three alternative methods for imputation: (1) regression method, (2) predictive mean matching, and (3) propensity score method. With binary (or ordinal) responses, the illustrative SAS commands in [Table 18.9](#) can be used to create 25 completed data sets using logistic regression imputations. Below we present a brief description of the command statements used in [Tables 18.8](#) and [18.9](#). For a more detailed description of these command statements, and other statements and options, the reader is referred to the extensive SAS documentation for PROC MI.

**Table 18.8** Illustrative commands for multiple imputation via MCMC method using PROC MI in SAS.

---

```
PROC MI DATA=widefile SEED=364865 NIMPUTE=25 OUT=mifile;
```

```
VAR group y1 y2 y3 y4;  
MCMC NBITER=5000 NITER=500;
```

---

**Table 18.9** Illustrative commands for multiple imputation via logistic regression using PROC MI in SAS.

```
PROC MI DATA=widefile SEED=364865 NIMPUTE=25 OUT=mofile;  
  VAR group y1 y2 y3 y4;  
  CLASS group y1 y2 y3 y4;  
  MONOTONE LOGISTIC(y2=group y1 group*y1);  
  MONOTONE LOGISTIC(y3=group y1 y2 group*y1 group*y2);  
  MONOTONE LOGISTIC(y4=group y1 y2 y3 group*y1 group*y2 group*y3);
```

---

PROC MI <options>;

The PROC MI statement calls the procedure MI in SAS. It includes options for the input SAS data-set to be read in (DALTA=SAS-data-set) and for the creation of an output SAS data-set (OUT=SAS-data-set) that contains the completed data sets with imputed values. The output SAS data-set includes an additional index variable, \_IMPUTATION\_, to identify the imputation number. For each imputation, the output data set contains all the variables in the input data set, with missing values replaced by imputed values. The PROC MI statement also includes options for the number of imputations to be created (NIMPUTE=*number*) and a seed used to initialize the random number generator (SEED=*number*); the latter is useful to ensure that results can be duplicated in a later run.

VAR variables;

The VAR statement lists the variables to be used in the imputation process. When the patterns of missingness are monotone, the order of variables in the VAR statement is important; they should be listed in order of the variables that are fully observed, followed by the variable with the fewest missing values, then the variable with the second fewest missing values, and so on. Note that the VAR statement can include variables that are fully observed but thought to be predictive of the probability of missingness and/or predictive of the missing responses.

CLASS variables;

The CLASS statement is used to define categorical variables in the VAR statement; these categorical variable can be used either as predictors for imputed variables or as imputed variables for data sets with monotone missing patterns. The CLASS statement must be used in conjunction with the MONOTONE statement (described later).

MCMC <options>;

The MCMC statement uses a Markov chain Monte Carlo method to impute values for a data set with an arbitrary missingness pattern, under an assumed multivariate normal distribution for the variables. It is also the default method of imputation when neither the MCMC nor the MONOTONE statement is specified. It includes options for controlling the number of burn-in iterations before the first imputation in each chain (MBITER=*number*) and the number of iterations between imputations in a chain (NITER=*number*).

MONOTONE <options>;

The MONOTONE statement specifies that the missingness patterns are monotone and provides three methods for imputing a quantitative variable: (1) regression method (REGRESSION), (2) predictive mean matching (REGPMM), and (3) propensity score method (PROPENSITY). These three options for the MONOTONE statement are:

REGRESSION <imputed = effects>

REGPMM <imputed = effects>

PROPENSITY <imputed = effects>

With the MONOTONE statement, the variables are imputed sequentially in the order given by the VAR statement. The effects specification option allows you to use a different set of predictors for each imputed variable (see [Table 18.9](#)). When the effect specification option is not used, the preceding variables in the VAR statement are used as the default predictors when imputing missing values for any particular variable. You can also specify more than one method of

imputation in the MONOTONE statement, and for each imputed variable, the predictors can be specified separately using the effects specification option.

The MONOTONE statement provides the following options for imputing a categorical variable (defined on the CLASS statement): (1) logistic regression (LOGISTIC) imputation of a binary or ordinal variable and (2) discriminant analysis (DISCRIM) of a nominal categorical variable. The effects specification option can also be used with both of these methods of imputation. The discriminant analysis imputation is based on assuming that within levels of the categorical variable the predictors have a multivariate normal distribution with means that vary across the categories (but with a constant covariance matrix); therefore it requires that all predictors of the missing categorical variable be continuous. In contrast, the predictors in a logistic regression imputation can be a mixture of continuous and categorical variables.

Finally, once the missing data are filled in  $m$  times, the  $m$  completed data sets can be analyzed using standard statistical methods. For example, the 25 completed data sets created in [Table 18.8](#) can be analyzed by fitting linear models using PROC MIXED in SAS. The results from these multiple analyses of the completed data sets can then be appropriately combined. [Table 18.10](#) presents illustrative commands for combining the results of analyses of multiple imputed data sets using PROC MIANALYZE in SAS. The first set of SAS commands in [Table 18.10](#) are used to convert the imputed data set (here named `mifile`) to a “long format” data set (`milong`) required for longitudinal analyses. Prior to analysis, the “long format” data set should be sorted by the index variable `_IMPUTATION_`. The next set of SAS commands use PROC MIXED to fit a linear mixed effects model, with randomly varying intercepts and slopes, to each of the imputed data sets. The BY statement in PROC MIXED is used to generate separate analyses based on observations grouped by the index variable `_IMPUTATION_`; this produces a separate analysis for each imputed data set. The ODS OUTPUT statement is used to create SAS data-sets containing the regression parameter estimates (`beta`) and the covariance matrices of the regression parameter estimates (`varbeta`) from the analyses of the imputed data sets. The final set of SAS commands use PROC MIANALYZE to appropriately combine these estimates to yield a single estimate of the regression parameters, together with standard errors that reflect the uncertainty inherent in the imputation of the missing data. The PARMS statement in PROC MIANALYZE is used to name the SAS data-set containing the regression parameter estimates (and the associated standard errors) from the analyses of the imputed data sets (here named `beta`); for multivariate inferences (e.g., multivariate Wald tests), the SAS data-set containing the covariance matrices of the regression parameter estimates (`varbeta`) must also be provided. For a more detailed description and explanation of command statements, and other options, the reader is referred to the extensive SAS documentation for PROC MIANALYZE.

**[Table 18.10](#)** Illustrative commands for combining the results of analyses of multiple imputed data using PROC MIANALYZE in SAS.

---

```
DATA milong;
SET mifile;
y=y1; time=0; OUTPUT;
y=y2; time=1; OUTPUT;
y=y3; time=2; OUTPUT;
y=y4; time=3; OUTPUT;
PROC SORT;
BY _IMPUTATION_;
PROC MIXED DATA=milong;
CLASS id;
MODEL y = group time group*time / S COVB;
RANDOM INTERCEPT time / SUBJECT=id TYPE=UN;
BY _IMPUTATION_;
ODS OUTPUT SOLUTIONF=beta COVB=varbeta;
PROC MIANALYZE PARMS=beta;
MODELEFFECTS INTERCEPT group time group*time;
```

---

# 18.7 COMPUTING: INVERSE PROBABILITY WEIGHTED (IPW) METHODS IN SAS

In principle, the application of inverse probability weighted (IPW) methods is straightforward because most statistical procedures allow for the inclusion of sampling weights. Specifically, the method can be implemented as a two-step process. The first step requires the calculation of the inverse probability weights by fitting a suitable model (e.g., logistic regression) for the probability of remaining in the study at each occasion (or, conversely, the probability of dropout). The second step requires the fitting of standard procedures for longitudinal analyses, where each observation is weighted by the inverse probabilities estimated in the first step. In the following we provide sample SAS commands for estimating the inverse probability weights and for replicating the weighted GEE analysis of the amenorrhea data from the *Clinical Trial of Contracepting Women* described in Section 18.4. The data for three subjects from the *Clinical Trial of Contracepting Women* are displayed in [Table 18.11](#).

**Table 18.11** Data for three subjects from the Clinical Trial of Contracepting Women.

<b>id</b>	<b>dose</b>	<b>time</b>	<b>y</b>	<b>prevy</b>	<b>r</b>
21	0	0	0	.	1
21	0	1	.	0	0
329	1	0	1	.	1
329	1	1	0	1	1
329	1	2	.	0	0
962	0	0	0	.	1
962	0	1	1	0	1
962	0	2	1	1	1
962	0	3	1	1	1

The first set of SAS commands in [Table 18.12](#) are used to read a “long format” data set, with separate records for each measurement occasion. The data set contains a variable, r, with  $r = 1$  if the outcome (y) is observed, and  $r = 0$  if missing due to dropout. In this illustration there are no missing data at baseline ( $time = 0$ ). The data set also contains a variable, prevy, denoting the value of the outcome at the previous occasion; by definition, prevy is missing at baseline ( $time = 0$ ) but is observed at all post-baseline occasions up to and including the time of dropout. The latter variable is used in the model for the inverse probability weights.

**Table 18.12** Illustrative commands for calculating cumulative probabilities of remaining in the study based on a logistic regression using PROC GENMOD in SAS.

```

DATA contracep;
  INFILE 'contracep.dat';
  INPUT id dose time y prevy r;
PROC SORT DATA=contracep;
  BY id time;
PROC GENMOD DESCENDING;
  CLASS time (PARAM=REF REF="1");
  MODEL r = time dose prevy dose*prevy / DIST=BIN;
  WHERE time NE 0;
  OUTPUT OUT=predict P=probs;
PROC SORT DATA=predict;
  BY id time;
DATA wgt (KEEP=id time cumprobs probs);
  SET predict;

```

```
BY id time;  
RETAIN cumprobs;  
IF FIRST.id then cumprobs=probs;  
ELSE cumprobs=cumprobs*probs;
```

---

The next set of SAS commands in [Table 18.12](#) use PROC GENMOD to fit a logistic regression model to the probability of remaining in the study at each occasion. The logistic regression is based on a “stacked” data set in which each individual contributes a sequence of binary “outcomes” ( $r$ ) to the analysis. Each binary “outcome”, denoted by  $R_{ij}$  in earlier sections of this chapter, is an indicator of whether the response was observed at a given occasion, from the second occasion (because, in this illustration, there are no missing data at baseline and  $R_{i1} = 1$  for all subjects) until either the occasion when dropout occurs or the last intended measurement occasion. The WHERE statement in PROC GENMOD restricts the analysis to the post-baseline occasions; the DESCENDING option ensures that the logistic regression models the probability that  $r = 1$ . The OUTPUT statement creates a new SAS data set named predict containing the estimated probabilities (probs) at each occasion. It is important that both the original SAS data-set (here named `contracep`) and the SAS data-set with the estimated probabilities (here named `predict`) are sorted by subject identification number (`id`), and within subject `id`, by measurement occasion (`time`).

The final set of SAS commands in [Table 18.12](#) are used to calculate the *cumulative* probabilities of remaining in the study at each occasion. This requires relatively sophisticated use of a DATA step, with both BY and RETAIN statements. In general, the use of the BY statement in a DATA step results in the creation of two temporary (so-called automatic) variables: FIRST.variable and LAST.variable (for any *variable* listed on the BY statement). In [Table 18.12](#), FIRST.`id` takes the value 1 for the first observation within `id` and 0 for all other observations within `id` (similarly LAST.`id` takes the value 1 for the last observation within `id` and 0 for all other observations within `id`). The RETAIN statement is used here to “remember” values of a variable from a previous observation. This use of the RETAIN statement allows for simple calculation of the cumulative probabilities (`cumprobs`) at each occasion for every subject.

Finally, the set of SAS commands in [Table 18.13](#) are used to: (1) merge the SAS data-set (`wgt`) containing the cumulative probabilities with the original SAS data-set (`contracep`), (2) create inverse probability weights, `ipw` (and set `ipw = 1` at baseline because there were no missing data at baseline), and (3) use PROC GENMOD to fit a marginal logistic regression model using IPW-GEE. The WEIGHT statement in PROC GENMOD weights each observation by the estimated inverse probability weights, `ipw`. To ensure that the weights are correctly applied within PROC GENMOD, a “working independence” assumption (TYPE=IND) for the within-subject association among the responses must be used. Because a “working independence” assumption is adopted, standard errors are based on the “sandwich” variance estimator.

**Table 18.13** Illustrative commands for calculating inverse probability weights and for fitting a marginal logistic regression model via IPW-GEE using PROC GENMOD in SAS.

```
DATA combine;  
MERGE contracep wgt;  
BY id time;  
IF (time=0) THEN ipw=1;  
ELSE ipw=1/cumprobs;  
PROC GENMOD DESCENDING DATA=combine;  
WEIGHT ipw;  
CLASS id;  
MODEL y = dose time time*time dose*time*dose / DIST=BIN;  
REPEATED SUBJECT=id / TYPE=IND;
```

---

## 18.8 FURTHER READING

General reviews of multiple imputation can be found in Rubin (1996), Rubin and Schenker (1991), Schafer (1999), and Horton and Lipsitz (2001). A comprehensive overview of the use of multiple imputation for handling incomplete data in longitudinal studies can be found in Chapters 9 and 13 of Molenberghs and Kenward (2007) and Chapter 28 of Molenberghs and Verbeke (2005). A detailed description of imputation by “chained equations” can be found in Raghunathan et al. (2001) and van Buuren (2007); also see van Buuren et al. (2006) for a study of the robustness of the method.

Inverse probability weighted methods have their roots in the survey sampling literature; IPW methods for longitudinal models were introduced in a landmark paper in the statistical literature by Robins et al. (1995). The article by Preisser, Lohman and Rathouz (2002) provides a concise but very useful summary of this topic in the context of dropout in longitudinal studies; also see Chapters 10 and 13 of Molenberghs and Kenward (2007).

Finally, missing data can arise due to death. Loss to follow-up due to death is qualitatively distinct from dropout due to other reasons and, ordinarily, needs to be handled quite differently in the analysis of longitudinal data; see Dufouil et al. (2004) for a very useful discussion of this topic.

# Bibliographic Notes

Multiple imputation was introduced by Rubin (1978) as a general method for handling missing data; see Rubin (1987) for a book-length treatment of this topic. Inverse probability weighted methods were first proposed in the sample survey literature by Horvitz and Thompson (1952). Robins et al. (1995) developed an inverse probability weighted estimating (IPW) equations approach for handling missing data in longitudinal studies. A more detailed description of the theory underlying the IPW methodology can be found in the text by Tsiatis (2006). The connection between imputation and inverse probability weighted methods is discussed in Reilly and Pepe (1997). Finally, Javaras and Fitzmaurice (2009) discuss analytic methods for handling extraneous covariates that are potentially predictive of missingness and also related to the covariates of interest and the outcomes.

<sup>1</sup> Implementations of the chained equation approach are available in the MICE library for R and S-Plus, ICE for Stata, and the IVEware macro for SAS (or standalone); a useful guide to software implementations is available at <http://multiple-imputation.com>.

<sup>†</sup> This section provides the formula for adjusting the “sandwich” variance estimator for estimation of the weights in IPW-GEE. The content of this section is somewhat technical and can be omitted without loss of continuity.

## *Part V*

# *Advanced Topics for Longitudinal and Clustered Data*

# *Chapter 19*

## *Smoothing Longitudinal Data: Semiparametric Regression Models*

### **19.1 INTRODUCTION**

The major focus of this book has been on parametric regression models for longitudinal data. Although parametric regression models (e.g., linear mixed effects models) have become established and enduring methods for longitudinal analyses, they have one important potential limitation: they assume that the shape of the functional relationship between the mean of the longitudinal response and the covariates is known. In general, these models describe covariate effects on the mean of the longitudinal response in terms of a relatively small number of regression parameters. This parsimony makes the description and interpretation of the results of a longitudinal analysis much simpler. However, in some settings the functional relationship between a covariate and the mean response may be too complex to be described by a few regression parameters. In such settings, what is required is a method of longitudinal analysis that allows greater flexibility for the form of the relationship. In nonparametric and semiparametric regression models the shape of the functional relationship between the response and covariate is not settled beforehand; instead, it is largely determined by the data at hand.

For example, in certain longitudinal studies it may be necessary to incorporate time trends in a completely nonparametric fashion, thereby allowing the mean response to change in a highly nonlinear, but not predetermined, way. In addition it may be appealing to assume a linear regression on other covariates; for example, treatment or exposure group contrasts might be modeled parametrically. Semiparametric regression models allow for such complexity by incorporating a component that is nonparametric (e.g., time trends) with a component that is parametric (e.g., treatment or exposure effects). Thus semiparametric regression models differ from nonparametric regression models in allowing the mean response (or a suitable known transformation of the mean response) to be modeled by nonparametric functions of certain covariates and parametric functions of other covariates.

In other settings nonparametric and semiparametric regression models may be used because they reduce reliance on *a priori* assumptions about the functional form of the relationship between the covariate and the response, while retaining the potential to suggest a suitable parametric model (e.g., linear or quadratic trend) that may fit the data at hand. That is, nonparametric and semiparametric regression methods can be quite useful in preliminary analyses, often guiding the appropriate choice of a parametric model that might fit the data well.

In this chapter we consider semiparametric regression techniques that do not require strong assumptions concerning the functional form of the pattern of change in the mean response. The methods that we describe can be broadly grouped together and referred to as “smoothing” methods. Although semiparametric regression models are well-developed for cross-sectional data with a single univariate response, the methodology has only recently been extended to the longitudinal setting. The correlation among repeated measures on the same individual poses additional complexity when extending these techniques to the longitudinal data setting.

The focus of much of this chapter is on a class of methods referred to as “penalized splines.” The motivation for this choice is twofold. First, penalized splines are relatively straightforward extensions of familiar linear regression models. Second, penalized splines are closely connected with linear mixed effects models. As we discuss later in this chapter, there is a mixed effects model

representation of penalized spline models that makes their extension to the longitudinal setting relatively straightforward and transparent. We begin this chapter by reviewing penalized splines for a single outcome in the cross-sectional setting. We then discuss extensions of these methods to longitudinal data.

## 19.2 PENALIZED SPLINES FOR A UNIVARIATE RESPONSE

In conventional parametric regression models the functional form of the relationship between the mean response and the covariates is assumed to be known; what is unknown is the values of the regression parameters. For example, the relationship between the mean response and a covariate might be assumed to be quadratic,

$$Y_i = \beta_1 + \beta_2 x_i + \beta_3 x_i^2 + e_i,$$

where  $Y_i$  and  $x_i$  denote the response and covariate for the  $i^{th}$  subject; the errors,  $e_i$ , are assumed to have mean zero. However, in some applications this is too restrictive, and a parametric regression model cannot faithfully capture the true underlying relationship between the response and covariate. In these settings nonparametric or semiparametric models are appealing because they allow the data to speak for themselves in determining the form of the relationship between the response and covariate. For example, nonparametric regression models only assume that the mean response is some smooth function of the covariate and allow the data to determine the form of the underlying function that relates these two variables.

In the univariate setting, nonparametric regression methods can be broadly divided into two main strands: kernel methods and spline-based methods. In this chapter we focus only on the latter, although we note in passing that there are some theoretical results that demonstrate close connections between kernel methods and certain spline methods. Specifically, we focus on smoothing methods known as *penalized splines*. To fix ideas, suppose that we have a response variable,  $Y_i$ , and a single covariate,  $x_i$  obtained on  $N$  individuals ( $i = 1, \dots, N$ ), and that it is of interest to understand the underlying relationship between  $Y_i$  and  $x_i$ . Consider the following simple nonparametric regression model:

$$(19.1) \quad Y_i = \theta(x_i) + e_i,$$

where  $\theta(x)$  is an unknown smooth regression function, describing the relationship between the mean of  $Y_i$  and the covariate,  $x_i$ . The errors,  $e_i$ , are assumed to be independent, with  $e_i \sim N(0, \sigma_e^2)$ . The goal is to estimate the nonparametric regression function,  $\theta(x)$ , from the data at hand. Although there is a plethora of nonparametric techniques for estimating the regression function  $\theta(x)$ , we focus exclusively on penalized spline estimation.

To motivate penalized splines, we begin by considering a smooth estimate of  $\theta(x)$  based on a piecewise linear function. Recall from Section 6.3 that a simple “broken-stick” model, with single knot at  $k$ , is given by,

$$Y_i = \beta_1 + \beta_2 x_i + \beta_3 (x_i - \kappa)_+ + e_i,$$

where  $(x_i - k)_+ = (x_i - k)$  if  $(x_i - k) > 0$  and is equal to zero otherwise. Naturally, a “broken-stick” model with only a single “break” or knot will not be sufficiently flexible to capture many non-linear relationships. A more general piecewise linear function is obtained by including additional knots, say  $k_1, \dots, k_M$ ,

$$Y_i = \beta_1 + \beta_2 x_i + \beta_{21} (x_i - \kappa_1)_+ + \beta_{22} (x_i - \kappa_2)_+ + \cdots + \beta_{2M} (x_i - \kappa_M)_+ + e_i.$$

The regression parameters for this piecewise linear model,  $\beta_1, \beta_2, \beta_{21}, \beta_{22}, \dots, \beta_{2M}$ , can be estimated via ordinary least squares (OLS) using standard statistical software for linear regression. Specifically, the OLS estimates are obtained by minimizing the residual sums of squares,

$$\sum_{i=1}^N \{Y_i - (\hat{\beta}_1 + \hat{\beta}_2 x_i + \hat{\beta}_{21} (x_i - \kappa_1)_+ + \hat{\beta}_{22} (x_i - \kappa_2)_+ + \cdots + \hat{\beta}_{2M} (x_i - \kappa_M)_+)\}^2.$$

The more knots that are included in the piecewise linear function, the greater is the flexibility of the fitted curve. Note that *locally* the approximation to the regression function,  $\theta(x)$ , has not increased in *complexity* in the sense that at any particular location, it is always a straight line regardless of the number of knots that have been included. However, a model with too many knots can yield a fitted

curve that may not be sufficiently smooth, that is, a fitted curve that may be too “rough” and variable. In addition, by including too many knots, there is a risk of overfitting the data. That is, from the fitted function it will be difficult to distinguish the true underlying trend from mere random fluctuations in the data.

In fitting flexible piecewise linear models, we face the challenge of choosing both the number and the location of the knots. Consequently it would be useful to have some way to choose among models with varying numbers of knots. Although standard selection criteria can be used to choose among models, these approaches quickly become prohibitive because the potential number of candidate models is very large. With  $M$  candidate knots, there are  $2^M$  possible models, such as with  $M = 20$  there are more than a million possible models to choose from. Fortunately, penalized spline regression provides a clever solution to this dilemma.

In penalized spline regression, we retain a relatively large number of knots in the model (although  $M$  is generally much smaller than  $N$ ) but *constrain* their influence. Instead of choosing which knots to retain by forcing some of the regression parameters,  $\beta_{21}, \beta_{22}, \dots, \beta_{2M}$ , to zero and allowing others to take on arbitrary values, we shrink all of the associated regression parameters toward zero. Specifically, we constrain the magnitudes of the regression parameters by requiring that

$$\sum_{m=1}^M \beta_{2m}^2$$

be less than some chosen positive value. Depending on the choice of value, this will lead to a smoother, less variable, fitted curve than would otherwise be obtained using the standard OLS criterion. To estimate the regression parameters subject to this constraint, we must minimize the following criterion:

$$\begin{aligned} & \sum_{i=1}^N \{Y_i - (\beta_1 + \beta_2 x_i + \beta_{21}(x_i - \kappa_1)_+ + \dots + \beta_{2M}(x_i - \kappa_M)_+)\}^2 \\ & + \lambda \sum_{m=1}^M \beta_{2m}^2, \end{aligned}$$

where  $\lambda > 0$ . This slight modification to the usual OLS criterion adds a “roughness penalty”,  $\lambda \sum_{m=1}^M \beta_{2m}^2$ . Note that only the coefficients for the *truncated line functions*  $(x_i - k_m)_+$ , the  $\beta_{2m}$ 's, are constrained, while the coefficients for the constant (1) and for  $x_i$  are left unpenalized. The “roughness penalty,”  $\lambda \sum_{m=1}^M \beta_{2m}^2$ , yields a smoother fit to the data depending on the magnitude of  $\lambda$ . For example, as  $\lambda \rightarrow 0$  (corresponding to infinite smoothing), the roughness penalty term predominates the criterion above and the impact of the knots diminishes. This results in a penalized spline estimator that effectively assumes  $\theta(x)$  is a simple linear function of  $x$ , and the penalized spline estimate converges to the OLS estimate of the linear regression of  $Y_i$  on  $x_i$ . Conversely, as  $\lambda \rightarrow 0$  (corresponding to no smoothing), the regression parameters are unconstrained and the penalized spline estimator is equivalent to the OLS estimator of the piecewise linear function, producing a rougher curve that potentially overfits the data. Because the amount of smoothing is determined by the value of  $\lambda$ , it is referred to as a *smoothing parameter*, in particular,  $\lambda$  governs the trade-off between smoothness and goodness-of-fit to the data. The smoothing parameter can be estimated from the data in a number of different ways.

Before discussing estimation of the smoothing parameter, we note that in the model above we have assumed a *linear* spline model for  $\theta(x)$ ,

$$(19.2) \quad \theta(x_i) = \beta_1 + \beta_2 x_i + \sum_{m=1}^M \beta_{2m}(x_i - \kappa_m)_+$$

Linear splines are simple but do not always provide a very smooth approximation to a complex function. Nonetheless, given a sufficient number of knots, linear splines will suffice in most applications. In principle, we can smooth out the corners at the knot locations by considering *polynomial* spline models of any order. For example, a *quadratic* spline model with knots at  $k_1, \dots,$

$k_M$  is given by

$$\theta(x_i) = \beta_1 + \beta_2 x_i + \beta_3 x_i^2 + \sum_{m=1}^M \beta_{3m} (x_i - \kappa_m)_+^2,$$

where  $(x_i - k_m)_+^2 = [(x_i - K_m)_+]^2$ . In general, a quadratic spline model will fit turning points (e.g., peaks and valleys) better than a linear spline model. However, with a sufficient number of knots, there is often little discernible difference between the fits provided by linear and quadratic spline models. Similarly a *cubic* spline model with knots at  $k_1, \dots, k_M$  is given by

$$\theta(x_i) = \beta_1 + \beta_2 x_i + \beta_3 x_i^2 + \beta_4 x_i^3 + \sum_{m=1}^M \beta_{4m} (x_i - \kappa_m)_+^3,$$

where  $(x_i - k_m)_+^3 = [(x_i - k_m)_+]^3$ . Cubic spline models are a common choice among higher-order polynomial splines. For ease of exposition, we focus only on *linear* spline models for the remainder of this chapter as they reveal most of the essential features of piecewise polynomial spline models yet are somewhat easier to understand and interpret.

Next we discuss estimation of the smoothing parameter  $\lambda$ . Interestingly, if the errors,  $e_i$ , are assumed to have a normal distribution, there is a close connection between the penalized spline estimator of  $\theta(x_i)$  and linear mixed effects models. Moreover this connection makes estimation of  $\lambda$  straightforward using existing statistical software. Specifically, the penalized spline estimator corresponds to a REML estimator in an equivalent linear mixed effects model. The linear mixed model representation of penalized splines results from expressing the regression function in terms of fixed and random effects,

$$(19.3) \quad \theta(x_i) = \beta_1 + \beta_2 x_i + \sum_{m=1}^M a_m (x_i - \kappa_m)_+,$$

where the random effects  $a_m \sim N(0, \sigma_a^2)$ . In this mixed model representation,

$$Y_i = \beta_1 + \beta_2 x_i + \sum_{m=1}^M a_m (x_i - \kappa_m)_+ + e_i,$$

the coefficients for the truncated line functions  $(x_i - k_m)_+$  are the random effects. Note that the model assumes only a *single* realization of  $(a_1, a_2, \dots, a_M)$ , and these  $M$  random coefficients are shared by all individuals. Moreover in this mixed model representation of penalized splines, it can be shown that

$$\lambda = \frac{\sigma_e^2}{\sigma_a^2}.$$

This representation of penalized splines is particularly useful because it allows estimation, and automatic determination of the degree of smoothing, to be placed firmly within a familiar mixed modeling framework that can be implemented using widely available statistical software. Specifically, it provides an automatic choice for the smoothing parameter via REML estimation of a variance parameter in the corresponding linear mixed effects model. In addition it suggests natural ways to extend penalized splines to the longitudinal setting via the incorporation of additional random effects in the model to account for the correlation among repeated measures.

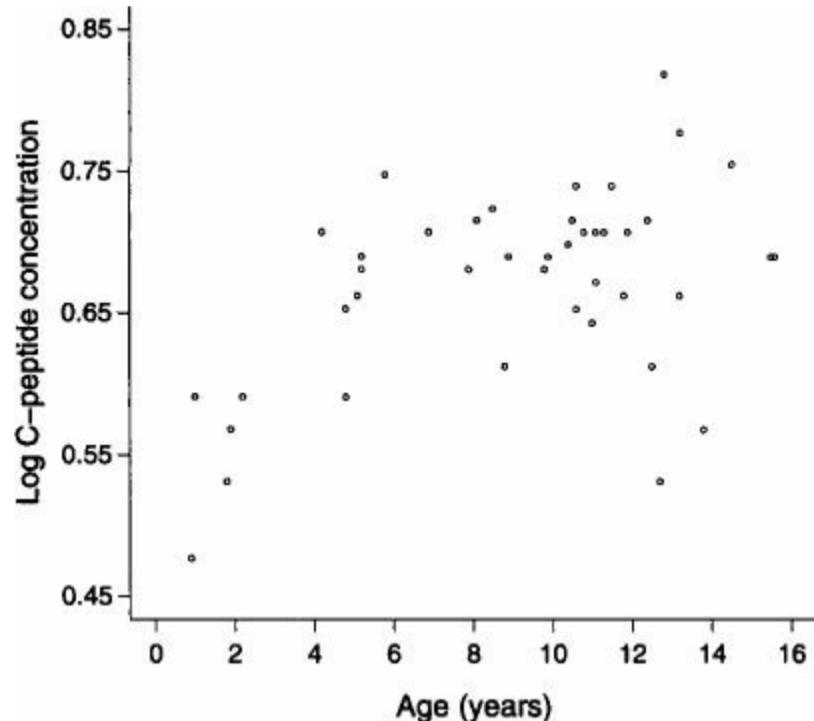
## 19.3 CASE STUDY

In this section we illustrate the use of penalized splines using cross-sectional data from a study of the factors that influence patterns of insulin-dependent diabetes mellitus in children (Sochett et al., 1987); the data are from Hastie and Tibshirani (1990, p. 304).

# Log C-Peptide Concentration in Children with Diabetes

In this study of 43 children, we are interested in the relationship between the logarithm of C-peptide concentration and age. [Figure 19.1](#) displays a scatterplot of log C-peptide concentration versus age and suggests that while log C-peptide concentration increases with age, the relationship appears to be non-linear. From [Figure 19.1](#) there is the suggestion that log C-peptide concentration increases until approximately age 7 or 8, but thereafter there may be a leveling off or even a decrease.

[Fig. 19.1](#) Scatterplot of log C-peptide concentration versus age in years for children with diabetes.



To better understand the relationship between log C-peptide concentration and age, we consider a penalized spline based on a piecewise linear curve with 10 knots at ages 5, 6, 7, 8, 9, 10, 11, 12, 13, and 14. Letting  $Y_i$  denote the log C-peptide concentration for the  $i^{th}$  child, the penalized spline is given by

$$E(Y_i|\text{Age}) = \beta_1 + \beta_2 \text{Age}_i + \sum_{m=1}^{10} \beta_{2m} \{\text{Age}_i - (m+4)\}_+,$$

with the magnitudes of the regression parameters constrained by the “roughness penalty,”  $\lambda \sum_{m=1}^{10} \beta_{2m}^2$ . Using a linear mixed effects model representation, this model can be expressed in an equivalent way as

$$E(Y_i|\text{Age}, a_1, \dots, a_{10}) = \beta_1 + \beta_2 \text{Age}_i + \sum_{m=1}^{10} a_m \{\text{Age}_i - (m+4)\}_+,$$

where the random effects  $a_m \sim N(0, \sigma_a^2)$ . Assuming that the errors,  $e_i = [Y_i - E(Y_i|\text{Age}, a_1, \dots, a_{10})]$ , are independent, with  $e_i \sim N(0, \sigma_e^2)$ , we can estimate the fixed effects and variance components ( $\sigma_e^2$  and  $\sigma_a^2$ ) by maximizing the likelihood. The REML estimates of the fixed effects and variance components are presented in [Table 19.1](#).

[Table 19.1](#) REML estimates of fixed effects and variance components (x 100) from mixed model representation of penalized spline for log C-peptide concentration in diabetic children.

Variable	Estimate	SE	Z
Intercept	0.5246	0.0317	16.53
Age	0.0269	0.0069	3.88
$\sigma_a^2 = \text{Var}(a_m)$	0.0104	0.0128	
$\sigma_e^2 = \text{Var}(e_i)$	0.3069	0.0699	

Based on the ratio of the estimates of the two variance components, the smoothing parameter is estimated to be  $\hat{\gamma} = \hat{\sigma}^2_e / \hat{\sigma}^2_a = 29.5$ . Given the estimates of the fixed effects and variance components, the BLUP estimates of  $a_m$  are presented in [Table 19.2](#). The BLUP estimates of  $a_m$  can be combined with the estimates of the fixed effects to yield the fitted curve. For example, the fitted mean at age 8 is given by,

$$\begin{aligned}\hat{Y} &= \hat{\beta}_1 + 8\hat{\beta}_2 + \hat{a}_1(8-5)_+ + \hat{a}_2(8-6)_+ + \cdots + \hat{a}_{10}(8-14)_+ \\ &= \hat{\beta}_1 + 8\hat{\beta}_2 + \hat{a}_1(8-5)_+ + \hat{a}_2(8-6)_+ + \hat{a}_3(8-7)_+ \\ &= 0.5246 + 0.0269 \times 8 - 0.0077 \times 3 - 0.0082 \times 2 - 0.0063 \times 1 \\ &= 0.694.\end{aligned}$$

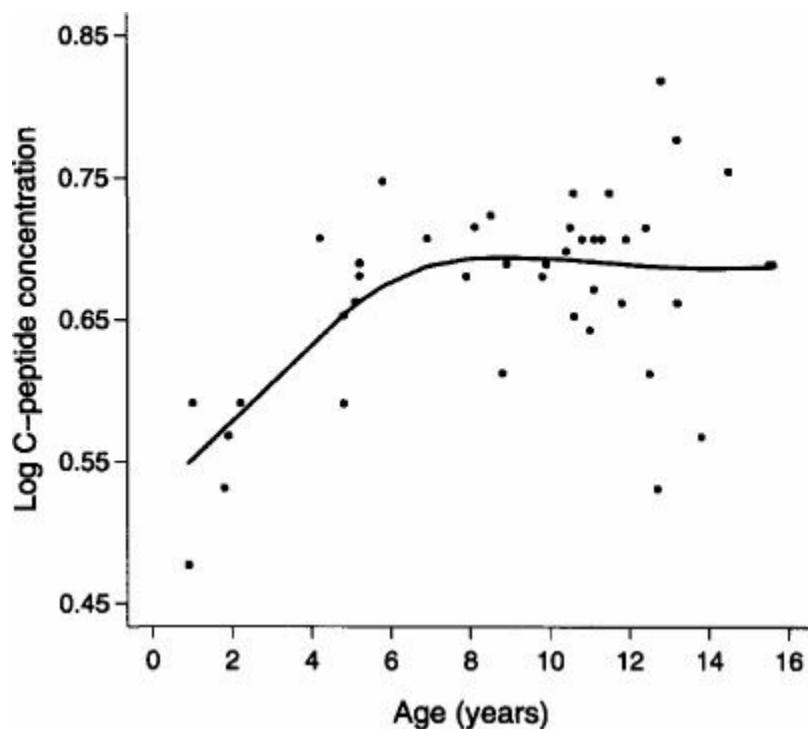
**Table 19.2** BLUP predictions of  $a_m$  ( $\times 100$ ) from mixed model representation of the penalized spline regression for log C-peptide concentration in diabetic children.

Variable	Estimate	SE of Prediction
$a_1$	-0.7728	0.8908
$a_2$	-0.8201	0.9147
$a_3$	-0.6290	0.9229
$a_4$	-0.3867	0.9086
$a_5$	-0.1574	0.9155
$a_6$	-0.1148	0.9112
$a_7$	-0.0463	0.9014
$a_8$	0.0616	0.9261
$a_9$	0.1050	0.9486
$a_{10}$	0.1376	0.9884

When the BLUP estimates of  $a_m$  are combined with the fixed effect estimates, we obtain the fitted curve presented in [Figure 19.2](#). It is apparent that log C-peptide concentration increases steadily until approximately age 8 and then levels off. Indeed, the fitted curve in [Figure 19.2](#) suggests that the relationship between log C-peptide concentration and age can be approximated by a piecewise quadratic-linear curve with single knot at age 8 (i.e., a quadratic trend prior to age 8, with linear trend thereafter). The corresponding regression model for this piecewise quadratic-linear curve is given by

$$E(Y_i | \text{Age}) = \beta_1 + \beta_2 \text{Age}_i + \beta_3 (\text{Age}_i - 8)_+ + \beta_4 \text{Age}_i^2 (8),$$

**Fig. 19.2** Time plot of log C-peptide concentration versus age in years, with fitted penalized spline superimposed, for children with diabetes.



where  $\text{Age}_i^2(8) = \text{Age}_i^2$  if  $\text{Age}_i \leq 8$  and  $\text{Age}_i^2(8) = 8^2 = 64$  if  $\text{Age}_i > 8$ ; that is,  $\text{Age}^2(8) = \min(\text{Age}^2, 8^2)$ . In [Table 19.3](#) we present the OLS estimates of the regression parameters from fitting this model to the diabetes data; the fitted curve, based on the regression parameter estimates, is displayed in [Figure 19.3](#). These results suggest that log C-peptide concentration increases significantly with age in the first 8 years but remains relatively constant thereafter. Specifically, the rate of change in log C-peptide concentration,  $(\hat{\beta}_2 + 2\hat{\beta}_4\text{Age})$ , is approximately 0.056 (or  $0.0658 - 2 \times 0.0047 \times 1$ ) at year 1, 0.019 (or  $0.0658 - 2 \times 0.0047 \times 5$ ) at year 5, and 0.0 (or  $0.0658 - 2 \times 0.0047 \times 7$ ) at year 7. Thus prior to age 8, C-peptide concentration is increasing at a rate of approximately 5.8% per year (or  $e^{0.056} = 1.058$ ) at age 1, 1.9% per year (or  $e^{0.019} = 1.019$ ) at age 5, and 0.0% per year (or  $e^{0.0} = 1.0$ ) at age 7. After age 8, C-peptide concentration remains relatively constant (with slope =  $0.0658 - 0.0663 = -0.0005$ ).

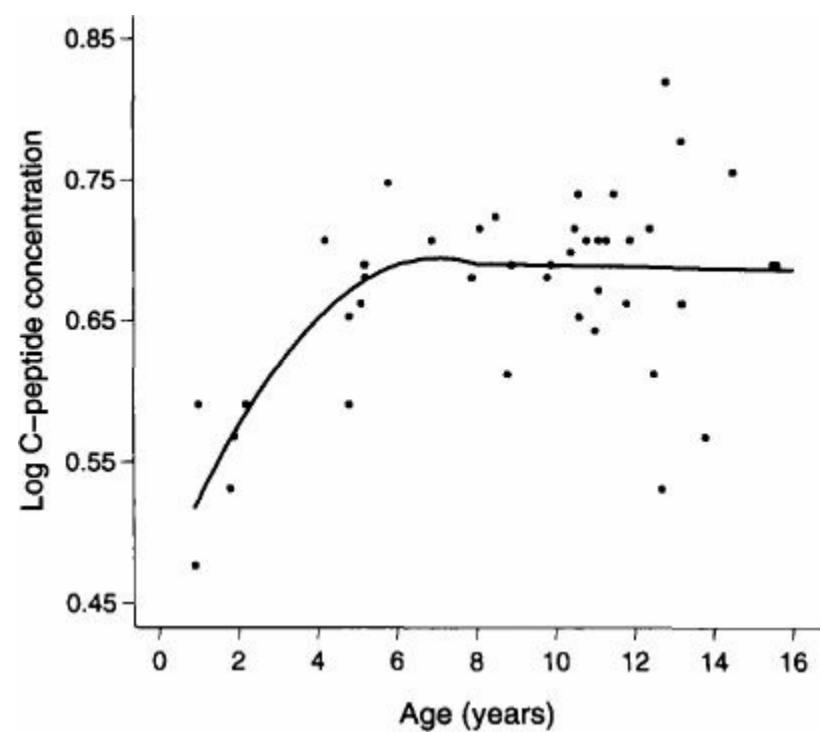
The fitted curve in [Figure 19.3](#), based on the regression parameter estimates in [Table 19.3](#), provides a fit to the data for ages 2 to 16 that is close to the fitted penalized spline in [Figure 19.2](#). The estimated rate of change in log C-peptide concentration from birth to age 5 should be very cautiously interpreted due to sparseness of data in that interval. That is, the fitted curve is based on few observations in that age interval (e.g., no observations between ages 2.2–4.2) and is therefore likely to be unreliable in that region.

**Table 19.3** OLS estimates and standard errors from piecewise quadratic-linear model, with single knot at age 8, for log C-peptide concentration in diabetic children.

Variable	Estimate	SE	Z
Intercept	0.4629	0.0456	10.15
Age	0.0658	0.0221	2.97
$(\text{Age} - 8)_+$	-0.0663	0.0214	-3.09
$\text{Age}^2(8)$	-0.0047	0.0023	-2.06

Note:  $\text{Age}^2(8) = \min(\text{Age}^2, 8^2)$

**Fig. 19.3** Time plot of log C-peptide concentration versus age in years, with fitted piecewise quadratic-linear curve superimposed, for children with diabetes.



## 19.4 PENALIZED SPLINES FOR LONGITUDINAL DATA

Next we consider smoothing methods for longitudinal data. A major challenge in extending smoothing techniques to longitudinal data is properly accounting for the within-subject correlation when constructing the penalized likelihood function. In a landmark book on semiparametric regression, Ruppert, Wand, and Carroll (2003) describe a very flexible semiparametric regression approach using the linear mixed model representation of penalized splines. In their approach the nonparametric component is kept relatively simple. This general approach can be extended in a relatively straightforward fashion to longitudinal data, while model fitting and inference can be embedded within an established and familiar linear mixed modeling framework.

Recall from Section 19.2 that the main idea underlying the linear mixed model representation of penalized splines for a *univariate* response is to express the mean as a function of fixed and random effects; the random effects are the coefficients for the truncated line functions. The linear mixed model representation of penalized splines for longitudinal responses expresses the mean response as a function of fixed and random effects in a similar way, but includes additional random subject effects to account for the correlation among the repeated measures. In these models the role of the random effects is two-fold, simultaneously handling two aspects of the data that require modeling: (1) one set of random effects, included in the model as coefficients for the truncated line functions, allows non-linear trends in the mean response, and (2) a second set of random effects, varying across subjects, accounts for the correlation among the repeated measures.

To fix ideas, consider the following piecewise linear function of time with  $M$  knots,  $k_1, \dots, k_M$ , and randomly varying subject effect (or random intercept):

$$(19.4) \quad Y_{ij} = \beta_1 + \beta_2 t_{ij} + \beta_{21}(t_{ij} - \kappa_1)_+ + \cdots + \beta_{2M}(t_{ij} - \kappa_M)_+ + b_i + \epsilon_{ij},$$

where  $Y_{ij}$  denotes the  $j^{th}$  response on the  $i^{th}$  individual at time  $t_{ij}$ ,  $b_i \sim N(0, \sigma^2_b)$  and  $\epsilon_{ij} \sim N(0, \sigma^2_\epsilon)$ . If we constrain the magnitudes of the regression parameters by adding a “roughness penalty,”  $\lambda \sum_{m=1}^M \beta_{2m}^2$ , the linear mixed model representation of penalized splines results from expressing the regression function above in terms of fixed and random effects:

$$(19.5) \quad \theta(t_{ij}) + b_i = \beta_1 + \beta_2 t_{ij} + \sum_{m=1}^M a_m(t_{ij} - \kappa_m)_+ + b_i,$$

where  $a_m \sim N(0, \sigma^2_a)$  and  $b_i \sim N(0, \sigma^2_b)$ . In this mixed model representation, the random effects,  $b_i$  (for  $i = 1, \dots, N$ ) and  $a_m$  (for  $m = 1, \dots, M$ ), have somewhat different roles. The inclusion of  $b_i$  allows each individual to have her own piecewise linear curve that is offset from the smooth population averaged curve by  $b_i$ ; marginally, this induces correlation among the repeated measures (albeit with the restrictive compound symmetry pattern). The additional random effects,  $a_m$ , are the coefficients for the truncated line functions,  $(t_{ij} - \kappa_m)_+$ , and produce a smooth regression function,  $\theta(t_{ij})$ , with the amount of smoothing depending on the relative magnitude of  $\sigma^2_a$ . Note that the two sets of random effects have different indices. There are distinct indices for  $a_m$  and  $b_i$  because the model assumes only a *single* realization of  $(a_1, a_2, \dots, a_M)$ . and these  $M$  random coefficients are shared by all individuals, whereas each individual is assumed to have a different random coefficient  $b_i$ . Estimation, and determination of the degree of smoothing, can be placed firmly within the familiar mixed modeling framework. Specifically, the fixed effects ( $\beta$ ) and variance parameters ( $\sigma^2_b$ ,  $\sigma^2_a$ , and  $\sigma^2_\epsilon$ ) can be estimated via REML, while the random effects, both  $a_m$  and  $b_i$ , can be predicted using BLUP.

The model presented above is very simple. It makes a strong assumption about the covariance among the repeated measures (i.e., compound symmetry). However, the model can be readily extended to allow for more general patterns of covariance through the inclusion of additional random

effects in the model. For example, the addition of random slopes to the model for the regression function,

$$(19.6) \quad \theta(t_{ij}) + b_{1i} + b_{2i}t_{ij} = \beta_1 + \beta_2 t_{ij} + \sum_{m=1}^M a_m(t_{ij} - \kappa_m)_+ + b_{1i} + b_{2i}t_{ij},$$

allows for heterogeneous variances and correlations that depend on  $t_{ij}$ . In principle, and when there is a sufficient number of repeated measures on each individual, it is even possible to allow each individual to have her own subject-specific, nonparametric smooth curve; we return to this topic in Section 19.6.

By including additional covariates in the regression model for the mean of  $Y_{ij}$ , the linear mixed model representation allows great flexibility for the forms of the relationships between the mean response and the covariates. For example, in the two group setting (e.g., active treatment versus control, or exposed versus non-exposed) time trends can be incorporated in a nonparametric fashion, thereby allowing the mean response to change in a highly non-linear, but not predetermined, way. This can be very appealing in the clinical trials setting where a pre-specified analysis plan is required but the functional form of the mean time trend cannot be settled beforehand. The group effect (e.g., treatment or exposure effect) can be incorporated in the model in a parametric fashion, thereby allowing a relatively simple, but powerful, test of the group effect on changes in the mean response over time. Consider the following model that illustrates these ideas,

$$Y_{ij} = \beta_1 + \beta_2 t_{ij} + \beta_3 \text{Group}_i + \beta_4 \text{Group}_i \times t_{ij} \\ + \sum_{m=1}^M a_m(t_{ij} - \kappa_m)_+ + b_{1i} + b_{2i}t_{ij} + \epsilon_{ij},$$

where  $\text{Group} = 1$  for the active treatment (or exposure) and 0 otherwise. This model allows a very general spline curve for the reference group ( $\text{Group} = 0$ ),

$$E(Y_{ij} | \text{Group}_i = 0) = \beta_1 + \beta_2 t_{ij} + \sum_{m=1}^M a_m(t_{ij} - \kappa_m)_+,$$

but constrains the *differences* between the smooth curves for the two groups to be a simple linear function of time,

$$E(Y_{ij} | \text{Group}_i = 1) - E(Y_{ij} | \text{Group}_i = 0) = \beta_3 + \beta_4 t_{ij}.$$

That is,  $\beta_3$  is the baseline difference between the mean response in the two groups (when  $t = 0$ ), whereas  $\beta_4$  is the constant rate of change (over time) in the *differences* between the smooth curves for the mean response in the two groups. Thus the assumed *linearity* refers to the hypothesized pattern of *differences* between the group means as a function of time and not to the time trends for the two groups. The inclusion of the random effects,  $a_m$ , constrains the coefficients for  $(t_{ij} - \kappa_m)_+$  and the relative magnitude of the variability of  $a_m$  determines the smoothness of the fitted curves. The inclusion of randomly varying intercepts and slopes,  $b_{1i}$  and  $b_{2i}$ , allows for heterogeneous variances and correlations that depend on  $t_{ij}$ . Note that in this model the comparison of groups in terms of their patterns of change in the mean response over time is based on a single-degree-of-freedom test of  $H_0: \beta_4 = 0$ . This test will have good statistical power to detect monotonically increasing (or decreasing) differences between the smooth curves for the two groups over time. We illustrate these ideas in the next section where smooth curves are fit to longitudinal data on progesterone metabolite concentration from a study of early pregnancy loss in two groups of women, a contraceptive and a non-conceptive group.

Finally, we conclude this section by noting that the focus of this chapter has been on smoothing methods for longitudinal data with a continuous response variable. The main ideas can be extended in a natural way to generalized linear models for discrete longitudinal data, allowing for smooth relationships in regression models for binary, ordinal, and count data. For example, in principle, it is relatively straightforward to develop a generalized linear mixed model (GLMM) representation of penalized splines for discrete data where the random effects include both random subject effects and random coefficients for the truncated line functions. However, in practice, these models are computationally more challenging to fit. With the inclusion of additional random coefficients for the

truncated line functions, the dimension of the vector of random effects is too large to rely on numerical quadrature for estimation and inference. To circumvent these computational challenges, estimation and inference are usually based on approximate methods (e.g., penalized quasi-likelihood (PQL); however, recall from Chapter 15 that PQL often yields badly biased estimates) or by using a Bayesian formulation of the GLMM fitted by Markov chain Monte Carlo (MCMC) algorithms.

## 19.5 CASE STUDY

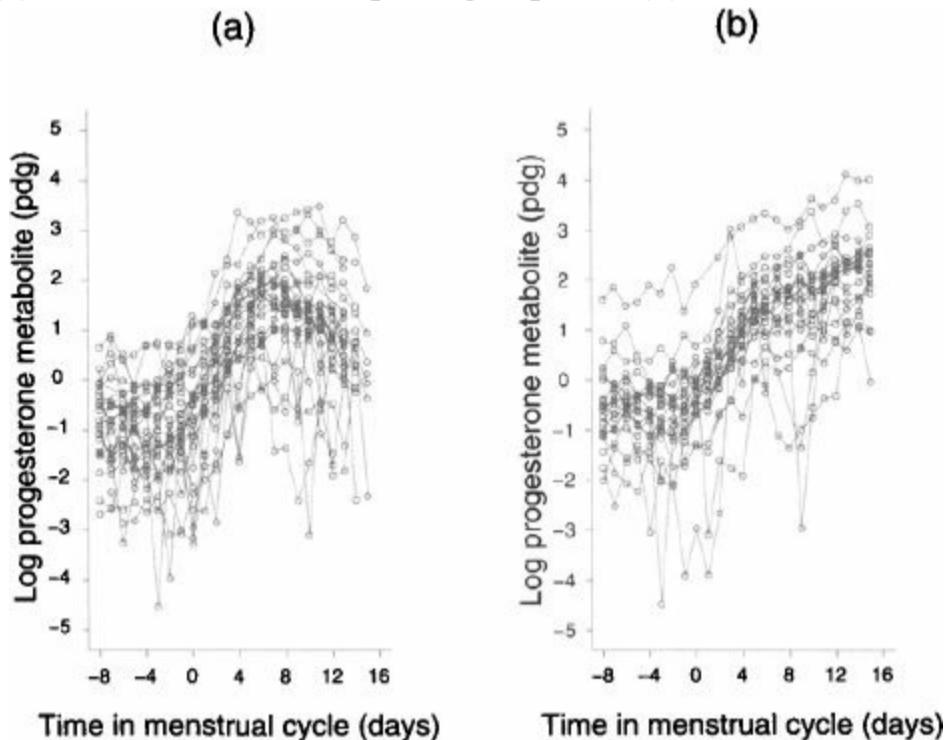
In this section we apply penalized regression splines to longitudinal data on a particular hormone, pregnanediol-3-glucuronide, measured daily in urine over the course of a menstrual cycle.

# Progesterone Metabolite Concentration

The hormone data are from a study of early pregnancy loss and consist of repeated progesterone metabolite (pregnanediol-3-glucuronide, PdG) measures from day  $-8$  to day  $15$  in the menstrual cycle (day  $0$  denotes ovulation day) on a sample of  $22$  conception cycles from  $22$  women and  $29$  non-conception cycles from another  $29$  women. Following standard practice in endocrinological research, progesterone profiles are aligned by the day of ovulation, determined by serum luteinizing hormone, and then truncated at each end to yield curves of equal length. Measures are missing for certain cycles on some days. The data are described in greater detail in Brumback and Rice (1998). There are additional data available on multiple non-conception cycles of the same  $29$  women, but for simplicity we focus on a single cycle. The goal of our analysis is to describe and compare the mean hormone profiles in the conception and non-conception groups, especially before and after implantation (implantation usually occurs approximately  $7$  to  $8$  days after ovulation).

We begin our analysis by considering separate time plots of the log progesterone metabolite concentrations for the conception and non-conception groups. [Figure 19.4](#) suggests that the trajectory of log progesterone concentration differs for the two groups following implantation (at approximately day  $7$ ). Specifically, for those in the non-conception group, log progesterone concentration appears to decrease beyond day  $7$ , while for those in the conception group log progesterone concentration appears to continue to increase. Moreover the overall shapes of the trajectories in the two groups seem to be highly non-linear.

**Fig. 19.4** Time plots, with joined line segments, of log progesterone concentration versus days of menstrual cycle for (a) women in non-conception group, and (b) women in conception group.



Next we fit separate penalized splines to the log progesterone metabolite concentrations for the conception and non-conception groups. Using the mixed effects model representation, and with  $22$  knots located consecutively from days  $-7$  through  $14$ , the model can be expressed as

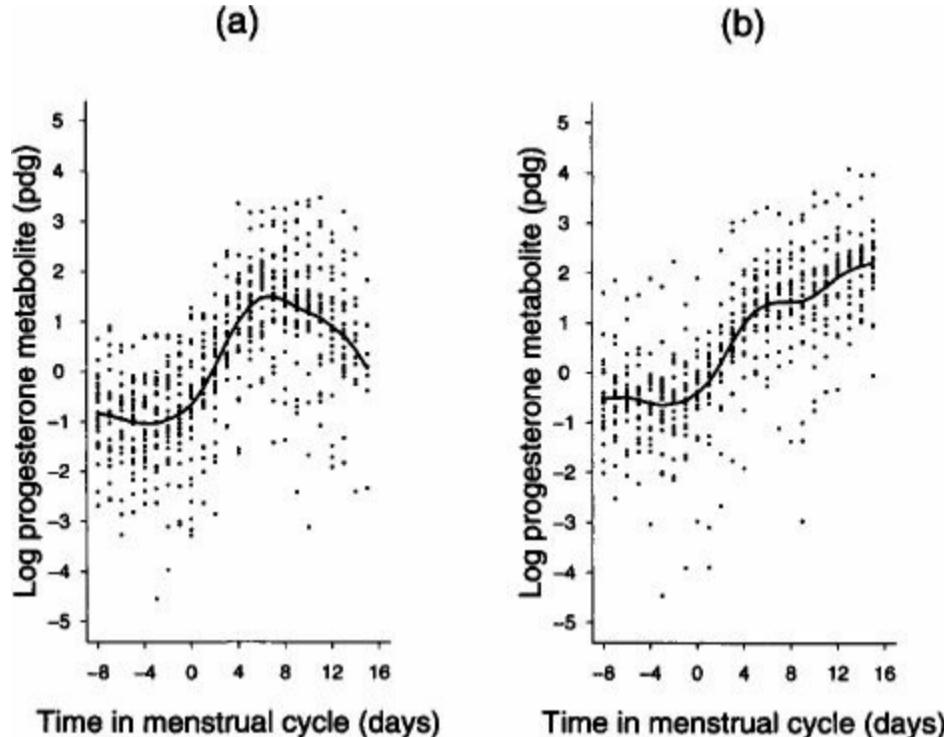
$$Y_{ij} = \beta_1 + \beta_2 t_{ij} + \sum_{m=1}^{22} a_m (t_{ij} - \kappa_m)_+ + b_{1i} + b_{2i} t_{ij} + \epsilon_{ij},$$

where  $Y_{ij}$  denotes the measure of log progesterone metabolite concentration for the  $i^{\text{th}}$  woman on day  $t_{ij}$  of her cycle. The inclusion of randomly varying intercepts and slopes,  $b_{1i}$  and  $b_{2i}$ , allows the variability of the response to be a quadratic function of time and also allows the magnitude of the correlation among repeated measures of log progesterone metabolite concentration to depend on the number of days separating them (see Section 8.3). The inclusion of the random effects,  $a_m$ , constrains the coefficients for  $(t_{ij} - k_m)_+$ . The random effects are assumed to have normal distributions, with  $a_m \sim N(0, \sigma_a^2)$  and

$$\begin{pmatrix} b_{1i} \\ b_{2i} \end{pmatrix} \sim N \left( \begin{bmatrix} 0 \\ 0 \end{bmatrix}, \begin{bmatrix} \sigma_{b1}^2 & \sigma_{b1,b2} \\ \sigma_{b1,b2} & \sigma_{b2}^2 \end{bmatrix} \right).$$

The fitted functions for the conceptional and non-conceptional groups, obtained by combining the REML estimates of the fixed effects and the BLUP predictions of the random effects,  $a_m$ , are displayed in [Figure 19.5](#). From a comparison of [Figure 19.5\(a\)](#) and [\(b\)](#), there are discernible differences between the two groups following implantation (at approximately day 7). For those in the non-conceptional group, the mean log progesterone concentration decreases beyond day 7, while for the conceptional group the mean log progesterone concentration increases. In addition, prior to day 7, the mean log progesterone concentration appears to be slightly higher in the conceptional group.

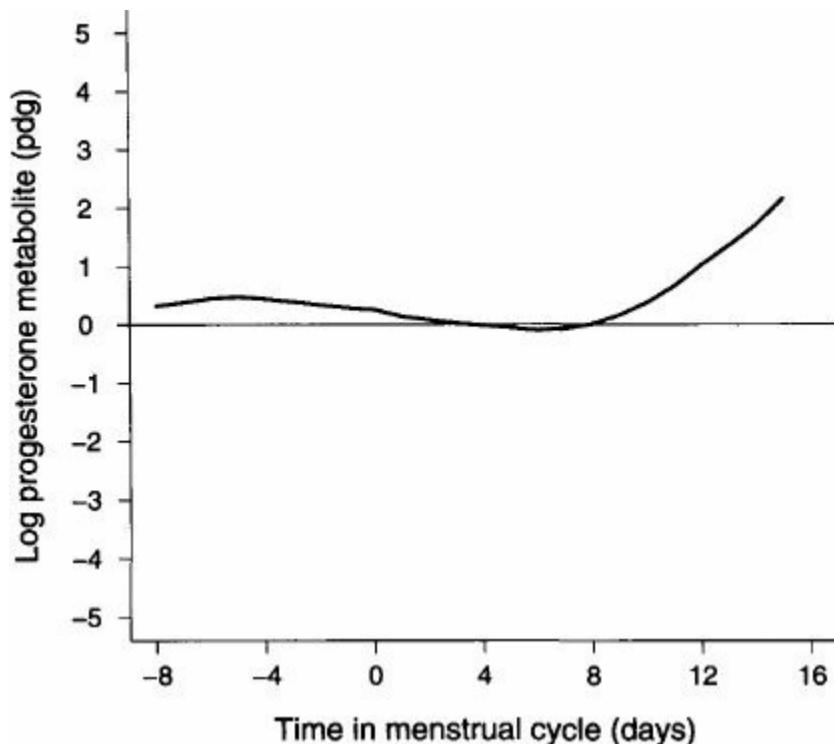
**Fig. 19.5** Time plots, with fitted penalized spline superimposed, of log progesterone concentration versus days of menstrual cycle for (a) women in non-conceptional group, and (b) women in conceptional group.



[Figure 19.6](#) displays the group differences (conceptional minus non-conceptional) in mean log progesterone concentration over time. This plot suggests large differences between the groups after implantation. The pattern in [Figure 19.6](#) also suggests that the most salient features of the *differences* between the smooth curves for the two groups can be represented by a piecewise linear trend with single knot at day 7, with a shallow decreasing slope prior to day 7, and a steep increasing slope after implantation on day 7. These features of the differences between the two groups can be captured in the following semiparametric regression model for the combined data from the conceptional and non-conceptional groups:

$$Y_{ij} = \beta_1 + \beta_2 t_{ij} + \beta_3 \text{Group}_i + \beta_4 \text{Group}_i \times t_{ij} + \beta_5 \text{Group}_i \times (t_{ij} - 7)_+ + \sum_{m=1}^{22} a_m (t_{ij} - \kappa_m)_+ + b_{1i} + b_{2i} t_{ij} + \epsilon_{ij},$$

**Fig. 19.6** Plot of group differences (conceptional versus non-conceptional group) in fitted penalized splines for log progesterone concentration during menstrual cycle.



where Group = 1 for the conceive group and 0 otherwise. This model allows a very general spline curve for the non-conceive group (the reference group) but constrains the differences between the curves for the conceive and non-conceive groups to be a simple piecewise-linear function of time, with different slopes before and after day 7. In this model the random effects are assumed to have normal distributions, with  $a_m \sim N(0, \sigma_a^2)$ ,

$$\begin{pmatrix} b_{1i} \\ b_{2i} \end{pmatrix} \sim N \left( \begin{bmatrix} 0 \\ 0 \end{bmatrix}, \begin{bmatrix} \sigma_{b_1}^2 & \sigma_{b_1, b_2} \\ \sigma_{b_1, b_2} & \sigma_{b_2}^2 \end{bmatrix} \right),$$

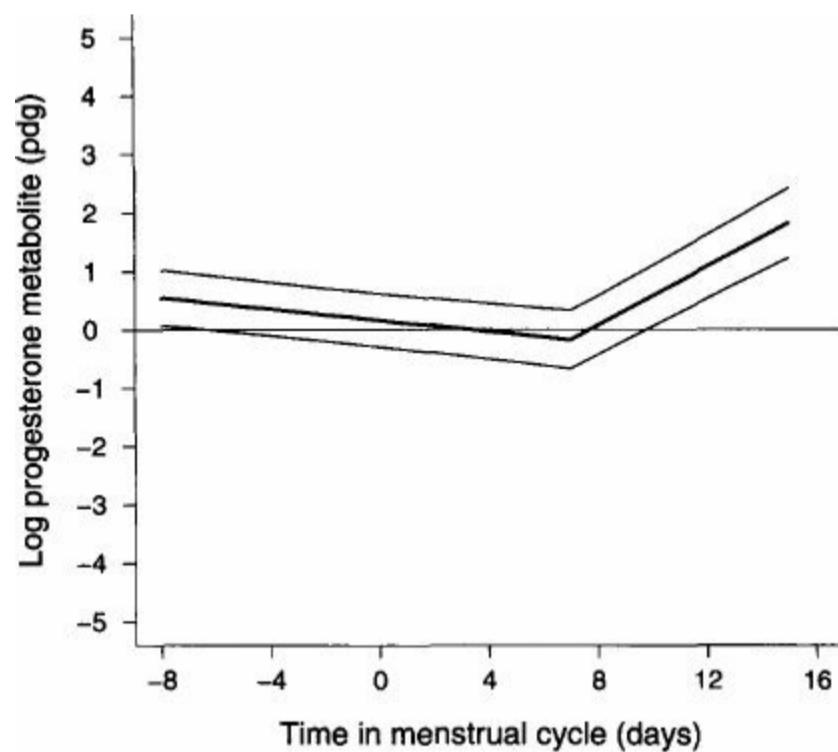
and the errors,  $\epsilon_{ij}$ , are assumed to be conditionally independent, with  $\epsilon_{ij} \sim N(0, \sigma_\epsilon^2)$ . Finally, this model implicitly assumes that the degree of smoothing is the same for both groups.

The REML estimates of the fixed effects and variance components are presented in [Table 19.4](#). These results indicate that the initial mean difference observed between the two groups in log progesterone concentration at the start of follow-up (day -8) has declined significantly prior to implantation (with estimated slope,  $\hat{\beta}_4 = -0.0479$ ,  $Z = -4.02$ ,  $p < 0.0001$ ). By day 7, there is no significant difference between the two groups in mean log progesterone concentration ( $\hat{\beta}_3 + 7 \times \hat{\beta}_4 = 0.1626 - 7 \times 0.0479 = -0.1727$ , SE = 0.2553,  $Z = -0.68$ ,  $p > 0.45$ ). Thereafter the trajectories of mean log progesterone concentration for the two groups depart significantly (with estimated slope,  $\hat{\beta}_4 + \hat{\beta}_5 = -0.0479 + 0.2962 = 0.2483$ , SE = 0.0203,  $Z = 12.2$ ,  $p < 0.0001$ ). [Figure 19.7](#) displays a plot of the estimated *differences* between the two groups (conceive versus non-conceive group), and 95% pointwise confidence limits.<sup>1</sup> It is transparent from this plot that the strongest evidence for differences between these two groups is during days 10 to 15, when the log progesterone concentration for the conceive group continues to increase while the log progesterone concentration for the non-conceive group decreases. Biologically, this result might be anticipated because hormones remain high after a successful conception or return to baseline levels when no conception occurs.

[Table 19.4](#) REML estimates and standard errors from mixed model representation of the semiparametric regression model for log progesterone concentration during menstrual cycle.

Variable	Estimate	SE	Z
Intercept	-0.8102	0.5713	-1.42
Time	0.0165	0.0752	0.22
Group	0.1626	0.2323	0.70
Group × Time	-0.0479	0.0119	-4.02
Group × (Time - 7) <sub>+</sub>	0.2962	0.0232	12.74
$\sigma_a^2 = \text{Var}(a_m)$	0.0162	0.0072	
$\sigma_{b1}^2 = \text{Var}(b_{1i})$	0.6576	0.1361	
$\sigma_{b2}^2 = \text{Var}(b_{2i})$	0.0010	0.0003	
$\sigma_{b1,b2} = \text{Cov}(b_{1i}, b_{2i})$	0.0038	0.0044	
$\sigma_\epsilon^2 = \text{Var}(\epsilon_{ij})$	0.2861	0.0127	

**Fig. 19.7** Plot of fitted group differences (non-conceptive versus conceptional group), and 95% confidence limits, for semiparametric regression model for log progesterone concentration during menstrual cycle.



# 19.6 FITTING SMOOTH CURVES TO INDIVIDUAL LONGITUDINAL DATA

So far the models we have considered have allowed the fitting of smooth curves to the *mean* of the longitudinal response. Next we consider extending these ideas to the fitting of smooth curves to each individual's data, thereby allowing the responses for each individual to be modeled as a subject-specific, non-linear function of time. Specifically, these smooth subject-specific curves can be modeled as penalized splines using the familiar mixed model representation described in previous sections.

Recall from Section 19.4 that the inclusion of randomly varying subject effects (e.g.,  $b_{1i}$  and  $b_{2i}$ ) allows each individual to deviate from the smooth population average curve. Marginally, when averaged over these random subject effects, this induces correlation among the repeated measures. Although the model allows each individual to deviate from the smooth population average curve, it makes the strong assumption that these deviations can be described parametrically. For example, in the context of the model given by (19.6) there is the assumption of linear departures,  $b_{1i} + b_{2i}t_{ij}$ , from the smooth population average curve.

To overcome this potential limitation we next consider an extension of the model that allows the subject-specific departures from the population average curve to be smooth functions, say  $f_i(t_{ij})$ . These functions are indexed by  $i$  to allow each individual to have her own subject-specific nonparametric function. To allow subject-specific departures from the population average curve to be smooth functions, the model given by (19.6) is extended as follows:

$$(19.7) \quad Y_{ij} = \theta(t_{ij}) + f_i(t_{ij}) + \epsilon_{ij},$$

with

$$\theta(t_{ij}) = \beta_1 + \beta_2 t_{ij} + \sum_{m=1}^M a_m (t_{ij} - \kappa_m)_+,$$

$$f_i(t_{ij}) = b_{1i} + b_{2i} t_{ij} + \sum_{k=1}^K u_{ki} (t_{ij} - \kappa_k)_+,$$

$$\text{where } a_m \sim N(0, \sigma_a^2), \quad \begin{pmatrix} b_{1i} \\ b_{2i} \end{pmatrix} \sim N\left(\begin{bmatrix} 0 \\ 0 \end{bmatrix}, \begin{bmatrix} \sigma_{b1}^2 & \sigma_{b1,b2} \\ \sigma_{b1,b2} & \sigma_{b2}^2 \end{bmatrix}\right),$$

$u_{ki} \sim N(0, \sigma_u^2)$ , and the errors,  $\epsilon_{ij}$ , are assumed to be conditionally independent, with  $\epsilon_{ij} \sim N(0, \sigma_\epsilon^2)$ . This model differs from (19.6) by extending the subject-specific departures from being simple linear functions of time, say  $b_{1i} + b_{2i}t_{ij}$ , to include a non-linear component of time,  $\sum_{k=1}^K u_{ki} (t_{ij} - \kappa_k)_+$ . In (19.7) the random effects  $a_m$ , included in the model as coefficients for the truncated line functions  $(t_{ij} - \kappa_m)_+$ , allow non-linear trends in the *population* mean response, whereas the random effects  $u_{ki}$ , included in the model as *subject-specific* coefficients for the truncated line functions  $(t_{ij} - \kappa_k)_+$ , allow subject-specific departures from the population average curve to be non-linear. In general, the number of truncated line functions specified for  $f_i(t_{ij})$  will need to be far fewer than are specified for  $\theta(t_{ij})$ ; that is, we require  $K \leq M$ . Because (19.7) can be placed in the standard linear mixed effect model framework, it is straightforward to include additional covariates in the regression model for the mean of  $Y_{ij}$ .

Finally, given the mixed effects model representation the fixed effects and variance components can be estimated using restricted maximum likelihood (REML). Given estimates of the fixed effects and the variance components, the random effects can be predicted using BLUR. When combined with the fixed effects, these predictions of the random effects yield smooth estimated curves for each individual and for the population averaged mean response.

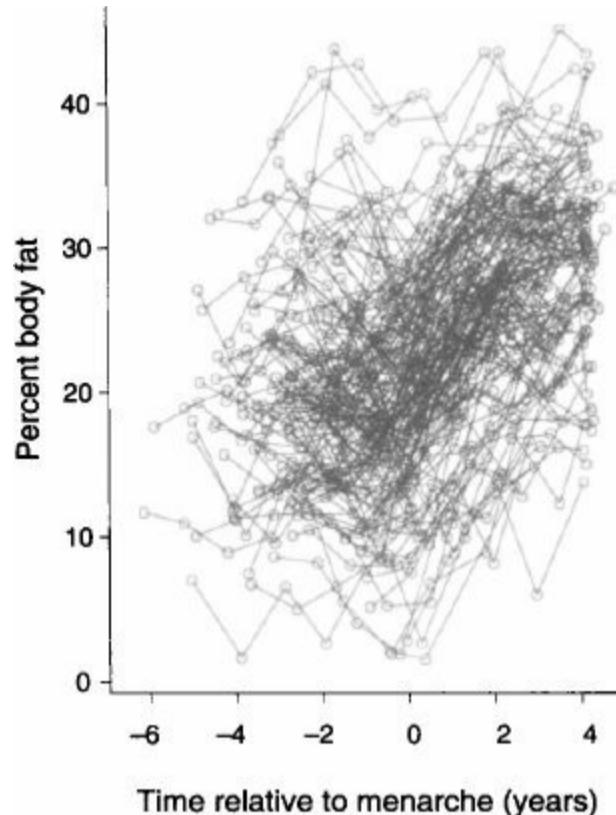
## 19.7 CASE STUDY

In this section we illustrate the fitting of smooth subject-specific curves to longitudinal data from the study on body fat accretion analyzed in Chapter 8, Section 8.8.

# Study of Influence of Menarche on Changes in Body Fat Accretion

Recall that these data are from a prospective study of the influence of menarche on body fat growth in a cohort of 162 girls from the MIT Growth and Development Study (Bandini et al., 2002; Phillips et al., 2003). [Figure 19.8](#) shows a time plot of the individual response profiles (where time is relative to the individual's age at menarche). From this graph it is difficult to discern whether the changes in percent body fat in the pre-menarcheal period are similar to the changes in the post-menarcheal period.

**Fig. 19.8** Time plot of percent body fat against time, relative to age of menarche (in years).



In Section 8.8 we presented an analysis of the changes in percent body fat before and after menarche. For the purposes of the analysis, “time” was coded as time since menarche. Specifically, we considered the following linear mixed effects model for the percent body fat data:

$$Y_{ij} = \beta_1 + \beta_2 t_{ij} + \beta_3(t_{ij})_+ + b_{1i} + b_{2i}t_{ij} + b_{3i}(t_{ij})_+ + \epsilon_{ij},$$

where  $t_{ij}$  denotes the time of the  $j^{th}$  measurement on the  $i^{th}$  subject before or after menarche (i.e.,  $t_{ij} = 0$  at menarche). In this model we assumed that each girl has a piecewise linear spline growth curve with a knot at the time of menarche. That is, each girl's growth curve can be described with an intercept and two slopes, one slope for the changes in response before menarche, another slope for the changes in response after menarche.

This analysis revealed a significant difference between the population slopes before and after menarche. Specifically, the estimate of the population mean pre-menarcheal slope was rather shallow, indicating that the annual rate of body fat accretion was less than 0.5%. In contrast, the estimate of the post-menarcheal slope was almost six times higher than the corresponding rate in the pre-menarcheal period and indicated that the annual rate of body fat accretion was approximately 2.5%. Finally, there was evidence of much heterogeneity among girls in their pre- and post-menarcheal slopes. Specifically, 95% of girls were predicted to have changes in percent body fat between -2.09% and 2.92% during the pre-menarcheal period, and between 0.62% and 4.30% during the post-menarcheal period.

Next we consider an analysis of these data that allows a smooth curve for the population mean response, in addition to smooth curves for each individual (i.e., smooth subject-specific curves). Specifically, we extend the model given above as follows:

$$Y_{ij} = \theta(t_{ij}) + f_i(t_{ij}) + \epsilon_{ij},$$

with

$$\theta(t_{ij}) = \beta_1 + \beta_2 t_{ij} + \beta_3(t_{ij})_+ + a_1(t_{ij} + 3)_+ + a_2(t_{ij} + 2)_+ + a_3(t_{ij} + 1)_+ + a_4(t_{ij} - 1)_+ + a_5(t_{ij} - 2)_+ + a_6(t_{ij} - 3)_+,$$

and

$$f_i(t_{ij}) = b_{1i} + b_{2i}t_{ij} + b_{3i}(t_{ij})_+ + u_{1i}(t_{ij} + 2)_+ + u_{2i}(t_{ij} + 1)_+ + u_{3i}(t_{ij} - 1)_+ + u_{4i}(t_{ij} - 2)_+,$$

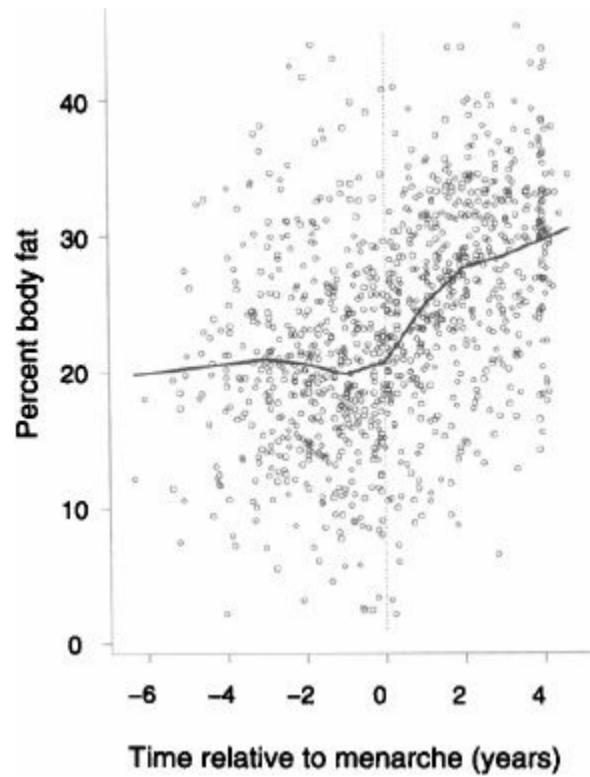
where  $a_m \sim N(0, \sigma^2_a)$  and  $u_{ki} \sim N(0, \sigma^2_u)$ . There are two important features of this model. First, it allows changes in the population average response to be non-linear functions of time both before and after menarche, via the inclusion of  $\sum_{m=1}^6 a_m(t_{ij} - \kappa_m)_+$ . From a subject-matter perspective, it might be anticipated that increases in the mean percent body fat begin to level off toward the end of adolescence. The model can accommodate this type of non-linearity in the post-menarcheal mean response. A second feature of this model is that it allows the subject-specific departures to be non-linear functions of time, both before and after menarche, via the inclusion of  $\sum_{k=1}^4 (t_{ij} - \kappa_k)_+$ . In this model we allow 7 truncated line functions (including the function  $(t_{ij})_+$  for the knot at menarche) for the smooth curve describing the population average mean response and 5 truncated line functions for describing the smooth, subject-specific departures from the population average curve.

The REML estimates of the fixed effects and variance components, and the BLUP predictions of  $a_1, \dots, a_6$  are presented in [Table 19.5](#). When the estimates of the fixed effects are combined with the BLUP predictions of  $a_1, \dots, a_6$ , they yield the fitted curve for the mean response displayed in [Figure 19.9](#). This fitted curve suggests that changes in the rate of growth in body fat occur 6 to 12 months prior to the onset of menarche and that the steepest rate of growth occurs in the 2 years after menarche. Of note, the estimate of  $\text{Var}(u_{ki})$ , 0.21 (with SE = 0.27), is relatively small when compared to the magnitudes of the other variance components. This suggests that subject-specific departures from the fitted mean curve can be approximated by a piecewise linear function with single knot at menarche (time = 0). For illustrative purposes, we display the predicted, smooth curves for two randomly selected individuals in [Figure 19.10](#).

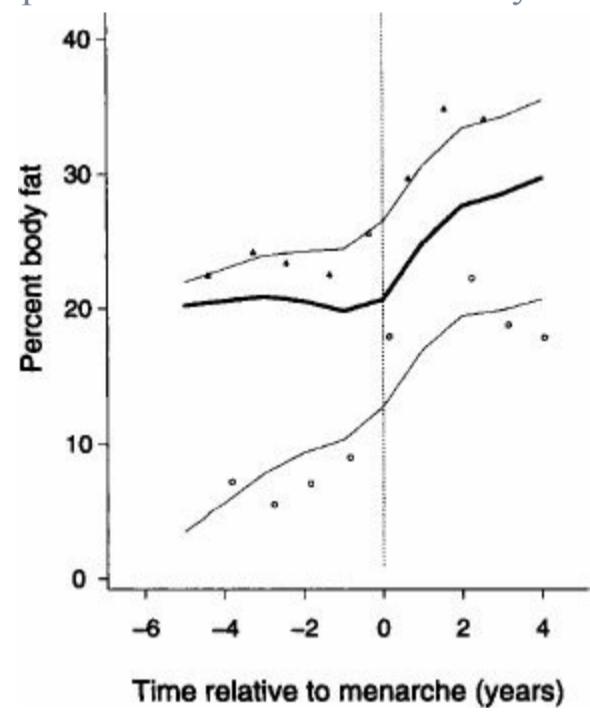
**Table 19.5** REML estimates and standard errors from mixed model representation of the semiparametric regression model with smooth subject-specific curves for percent body fat.

Variable	Estimate	SE	Z
Intercept	21.9412	1.5478	14.18
Time	0.3376	0.4462	0.76
(Time) <sub>+</sub>	3.3222	0.7515	4.42
$a_1$	-0.6826	0.6952	-0.98
$a_2$	-0.4112	0.7144	-0.58
$a_3$	1.6431	0.7039	2.33
$a_4$	-1.4277	0.7046	-2.03
$a_5$	-1.9478	0.6926	-2.81
$a_6$	0.3929	0.7554	0.52
$\text{Var}(a_m)$	2.0599	1.4423	
$\text{Var}(b_{1i})$	45.8487	5.7728	
$\text{Var}(b_{2i})$	1.8413	0.5085	
$\text{Var}(b_{3i})$	3.4590	1.0765	
$\text{Cov}(b_{1i}, b_{2i})$	2.5819	1.2919	
$\text{Cov}(b_{1i}, b_{3i})$	-6.9389	1.9378	
$\text{Cov}(b_{2i}, b_{3i})$	-2.4127	0.6601	
$\text{Var}(u_{ki})$	0.2132	0.2704	
$\text{Var}(\epsilon_{ij})$	8.0473	0.5051	

**Fig. 19.9** Time plot of percent body fat against time, relative to age of menarche (in years), with predicted smoothed curve.



**Fig. 19.10** Predicted smooth curve for mean percent body fat again time, relative to age of menarche (in years), with smooth subject-specific curves for two randomly selected subjects.



## 19.8 COMPUTING: FITTING SMOOTH CURVES USING PROC MIXED IN SAS

By exploiting their mixed effects model representation, the penalized splines discussed in Sections 19.2 through 19.5 can be fit using the PROC MIXED procedure in SAS. The linear mixed model representation of penalized splines expresses the regression function in terms of fixed and random effects. In the cross-sectional setting (see Section 19.2), the linear mixed model representation is

$$\theta(x_i) = \beta_1 + \beta_2 x_i + \sum_{m=1}^M a_m (x_i - \kappa_m)_+,$$

where the random effects are the coefficients for the truncated line functions  $(x_i - \kappa_m)_+$ . Note that the random effects are indexed by  $m$  ( $m = 1, \dots, M$ ), where  $M$  is the number of truncated line functions, and not by  $i$ , the index for subjects. The mixed effects model representation of penalized splines is completed by assuming the random effects  $a_m \sim N(0, \sigma_a^2)$ .

For example, the illustrative SAS commands given in [Table 19.6](#) can be used to fit a penalized spline to examine the relationship between the outcome (denoted by  $y$ ) and age in a cross-sectional study of children. The SAS commands in [Table 19.6](#) fit a penalized spline based on a piecewise-linear curve with 8 equally spaced knots at ages 5, 6, 7, 8, 9, 10, 11, and 12. The first part of the command syntax in [Table 19.6](#) is for the construction of the 8 truncated line functions. The truncated line functions for the linear spline model are denoted by  $bf1, bf2, \dots, bf8$ , where, for example,  $bf1 = (age - 5)_+$ . The second part of the command syntax is for the fitting of the linear mixed model representation of the penalized spline. The syntax for the RANDOM statement requires additional explanation. Because the random effects are the coefficients for the truncated line functions,  $bf1, bf2, \dots, bf8$  are included on the RANDOM statement. However, to ensure that the random effects  $a_1, a_2, \dots, a_8$  that multiply these functions are drawn from the same normal distribution, say  $a_m \sim N(0, \sigma_a^2)$ , we must constrain the covariance matrix of the 8 random coefficients to be  $\sigma_a^2 I_8$ , where  $I_8$  denotes an  $8 \times 8$  identity matrix. That is, the random coefficients must be constrained to have the same variance,  $\sigma_a^2$ . This is achieved by specifying TYPE=TOEP(1) for the structure of the covariance of the random effects. Recall from Chapter 7 that TOEP(1) denotes a banded Toeplitz matrix, with a band size of 1; this implies a covariance matrix with a constant variance along the diagonal and all covariances set to zero. Also note that the model assumes only a *single* realization of  $(a_1, a_2, \dots, a_8)$  and that these random coefficients are shared by all individuals; consequently the SUBJECT option on the RANDOM statement is not used. Recall that the SUBJECT option on the RANDOM statement determines when new realizations of the random effects are assumed to occur. Finally, the OUTPRED option on the model statement specifies an output data set (yhat) that contains, among other useful quantities, the BLUP predictions (PRED) of the mean response,  $\hat{y} = \hat{\theta}(x_i)$ . The BLUP predictions can be used to display the fitted smooth curve by plotting PRED against pre-sorted values of age; see the third part of the command syntax in [Table 19.6](#).

**Table 19.6** Illustrative commands for fitting a penalized spline to cross-sectional data using PROC MIXED in SAS.

---

```

bf1 = max(0, age-5);
bf2 = max(0, age-6);
bf3 = max(0, age-7);
bf4 = max(0, age-8);
bf5 = max(0, age-9);
bf6 = max(0, age-10);
bf7 = max(0, age-11);
bf8 = max(0, age-12);

PROC MIXED;
  MODEL y=age / SOLUTION OUTPRED=yhat;
  RANDOM bf1 bf2 bf3 bf4 bf5 bf6 bf7 bf8 / TYPE=TOEP(1) SOLUTION;

```

```

PROC SORT DATA=yhat;
  BY age;
SYMBOL INTERPOL=JOIN;

PROC GPLOT DATA=yhat;
  PLOT PRED*age;

```

---

In the longitudinal setting (see Section 19.4), an example of the linear mixed model representation for the regression function is

$$\theta(t_{ij}) + b_{1i} + b_{2i}t_{ij} = \beta_1 + \beta_2 t_{ij} + \sum_{m=1}^M a_m(t_{ij} - \kappa_m)_+ + b_{1i} + b_{2i}t_{ij},$$

where  $a_m \sim N(0, \sigma_a^2)$  and  $(b_{1i}, b_{2i})$  are assumed to have a zero mean bivariate normal distribution. Note that the two sets of random effects,  $a_m$  and  $(b_{1i}, b_{2i})$ , have distinct indices. The random effects  $a_m$  (for  $m = 1, \dots, M$ ) are the coefficients for the truncated line functions,  $(t_{ij} - \kappa_m)_+$ , and produce a smooth regression function  $\theta(t_{ij})$  that describes the population average response. The additional random effects  $(b_{1i}, b_{2i})$  allow each individual to deviate from the smooth population average curve; marginally, this induces a random effects covariance pattern among the repeated measures of the response (with heterogeneous variances and correlations that depend on  $t_{ij}$ ).

For example, the illustrative SAS commands given in [Table 19.7](#) can be used to fit a penalized spline to examine how the outcome (denoted by  $y$ ) changes over time in a longitudinal study. The SAS commands in [Table 19.7](#) fit a penalized spline with 8 knots at equally spaced times. The first part of the command syntax in [Table 19.7](#) is for the construction of the 8 truncated line functions. The second part is for the fitting of the penalized spline. Here two RANDOM statements are required because the two sets of random effects,  $a_m$  (for  $m = 1, \dots, 8$ ) and  $(b_{1i}, b_{2i})$  (for  $i = 1, \dots, N$ ) have different indices. The first RANDOM statement is for the specification of  $a_m$ . To ensure that the random effects  $a_1, a_2, \dots, a_8$  that multiply the truncated line functions are drawn from the same distribution, say  $a_m \sim N(0, \sigma_a^2)$ , we must constrain the covariance matrix of these random effects to be  $\sigma_a^2 I_8$ , where  $I_8$  denotes an  $8 \times 8$  identity matrix. This is achieved by specifying TYPE=TOEP(1) for the structure of their covariance matrix. Note that the model assumes only a single realization of  $(a_1, a_2, \dots, a_8)$ ; consequently the SUBJECT option on the RANDOM statement is not used. The second RANDOM statement is used to indicate that  $(b_{1i}, b_{2i})$  are drawn from a bivariate normal distribution with unstructured covariance matrix (TYPE=UN) and that  $(b_{1i}, b_{2i})$  vary across individuals according to the distinct values of id (SUBJECT = id); that is, each subject has a different pair of random coefficients,  $(b_{1i}, b_{2i})$ .

**Table 19.7** Illustrative commands for fitting a penalized spline to longitudinal data using PROC MIXED in SAS.

---

```

bf1 = max(0, time-1); bf2 = max(0, time-2);
bf3 = max(0, time-3);
bf4 = max(0, time-4);
bf5 = max(0, time-5);
bf6 = max(0, time-6);
bf7 = max(0, time-7);
bf8 = max(0, time-8);

PROC MIXED;
  CLASS id;
  MODEL y=time / SOLUTION;
  RANDOM bf1 bf2 bf3 bf4 bf5 bf6 bf7 bf8 / TYPE=TOEP(1) SOLUTION;
  RANDOM INTERCEPT time / SUBJECT=id TYPE=UN;

```

---

## 19.9 FURTHER READING

There is probably no better source for additional information on semiparametric regression, and the linear mixed model representation of penalized splines, than the excellent text by Ruppert, Wand, and Carroll (2003). Ruppert, Wand, and Carroll (2003) provide a truly expository guide to penalized splines, with numerous examples of how to fit the models using the mixed model representation. Gurrin, Scurrah, and Hazelton (2005) also provide a very accessible tutorial article on penalized spline smoothing using linear mixed models.

# Bibliographic Notes

The mixed model representation for fitting smooth curves has appeared in numerous sources in the statistical literature; see Ruppert, Wand, and Carroll (2009) for a comprehensive review of the state of the art in this field. Early contributions to the literature on the mixed model representation include Wahba (1978), Green (1985), Thompson (1985), Speed (1991), Donnelly, Laird, and Ware (1995), O'Connell and Wolfinger (1997), and Wang (1998). More recent developments focusing on longitudinal data include Welsh, Lin, and Carroll (2002), Carroll et al. (2004), Crainiceanu and Ruppert (2004), Ghidey, Lesaffre, and Eilers (2004), Durbán et al. (2005), Harezlak et al. (2005), Wang, Carroll, and Lin (2005), Welham et al. (2006), and Harezlak, Naumova, and Laird (2007).

<sup>1</sup> Each of these confidence intervals covers the corresponding true value of log progesterone metabolite concentration with probability 0.95. Taken together, these confidence intervals constitute a 95% *pointwise* confidence band for log progesterone metabolite concentration as a function of time. This is in contrast to a 95% *simultaneous* confidence band, where the probability is 0.95 that *all* of them cover their corresponding true values *simultaneously*; in general, the latter will be wider than a pointwise confidence band.

# *Chapter 20*

## *Sample Size and Power for Longitudinal Studies*

### **20.1 INTRODUCTION**

The emphasis in earlier chapters has been on methods for analyzing longitudinal data. In this chapter we consider the design of a longitudinal study. Specifically, we focus on the determination of sample size and power for longitudinal studies. In general, questions about sample size and power arise in the earliest stages of the design of a study. Although the question can be posed in a variety of different ways, investigators typically need to know the answer to the following question: How *large* should my study be? In a cross-sectional study design, with only a single univariate response, the answer to this question is relatively straightforward: the *size* of a study is directly related to the number of subjects, that is, the sample size. However, for a longitudinal study the question of *size* is more complex. For example, in planning a longitudinal study to compare an active treatment to control, investigators need to determine not only how many subjects to enroll in the study but also the duration of the study and the frequency and spacing of repeated measurements on the subjects.

For the special case of a cross-sectional study design, with only a single response, statisticians have developed simple formulas for sample size and power calculations. Explicit formulas can be found in many introductory textbooks in statistics. In addition some statistical software packages include procedures for sample size and power calculations, and publicly available sample size and power calculators can be found on the Web. However, for the multivariate response obtained from a longitudinal study, accurate sample size (and power) determination is more complicated and, in general, requires inversion of matrices and iterative solutions when no closed-form expressions can be obtained. The purpose of this chapter is not to derive complex sample size formulas for longitudinal studies; references to accurate, but also more complex, methods for calculating sample size and power can be found at the end of the chapter. Instead, we present simple, albeit approximate, methods for sample size and power determination for longitudinal studies that allow direct application of standard sample size and power formulas.

In this chapter we begin with a review of sample size (and power) formulas for a univariate continuous response in a cross-sectional study design. We emphasize the main considerations in determining how large the sample size needs to be to achieve a specified power to detect some effect of scientific interest; this section can be skimmed through for those already familiar with power and sample size calculations for a univariate response. We then present simple closed-form expressions for sample size (and power) calculations for longitudinal studies with a continuous response based on the standard sample size (and power) formula for a univariate response. Similar closed-form expressions for longitudinal binary responses are also presented. Finally, two examples are presented to illustrate the application of these formulas to the design of a longitudinal study with a continuous and binary response, respectively.

## 20.2 SAMPLE SIZE FOR A UNIVARIATE CONTINUOUS RESPONSE

When planning a cross-sectional study, investigators must establish how many subjects they will need to achieve some specified power to detect an effect of subject-matter importance. For example, suppose that investigators are interested in comparing two treatments, an active treatment and a control. The investigators plan to randomize a total of  $N$  subjects, with  $N_1 = \pi N$  in group 1 (e.g., active treatment), and  $N_2 = (1 - \pi)N$  in group 2 (e.g., control). When  $\pi = 0.5$ , an equal number of subjects ( $N_1 = N_2$ ) are randomized to receive each of the two treatments. At the completion of the study, the two treatment groups are to be compared in terms of the mean response. Let  $\mu^{(1)}$  denote the mean response in the population of individuals assigned to the active treatment; similarly let  $\mu^{(2)}$  denote the mean response in the population of individuals assigned to the control. The treatment effect can be expressed in a variety of different ways, but here we consider the simple difference in means,  $\delta = \mu^{(1)} - \mu^{(2)}$ . The null hypothesis of no treatment difference is represented by  $H_0: \delta = 0$ . In this example the investigators may be interested in establishing whether the active treatment is superior to control, with the alternative hypothesis that  $\delta > 0$ .

Before we discuss sample size and power, we must consider the two types of errors that can arise when conducting a statistical test of  $H_0: \delta = 0$ . The first kind of error is called a type I error and is made if we reject the null hypothesis when in fact it is true. The probability of a type I error, also known as the significance level of the test, is usually denoted by  $\alpha$ . Thus, for our example where  $H_0: \delta = 0$ ,

$$\alpha = \Pr(\text{Reject } H_0 \mid H_0 \text{ is true}).$$

Conventionally,  $\alpha$  is chosen to be no greater than 0.05; that is, we are prepared to mistakenly reject the null hypothesis no more than 5% of the time. The second kind of error that can arise when conducting a statistical test is called a type II error. A type II error is made if we fail to reject the null hypothesis when in fact it is false. We denote the probability of a type II error by  $\gamma$ , with

$$\gamma = \Pr(\text{Fail to reject } H_0 \mid H_0 \text{ is false}).$$

(The usual convention is to denote the probability of a type II error by  $\beta$ ; however, we have chosen to denote it by  $\gamma$  to avoid any potential confusion with our widespread use of  $\beta$  for the regression parameters in earlier chapters of the book.) Since  $\gamma$  is determined by considering the case where the null hypothesis is not true (i.e.,  $\delta \neq 0$ ), it necessarily depends on the particular choice of value for  $\delta \neq 0$  under the alternative hypothesis. Intuitively, the closer the true value of  $\delta$  is to zero (the assumed value for  $\delta$  under the null hypothesis), the more difficult it is to reject  $H_0: \delta = 0$ . Finally, the power of a statistical test is defined as  $1 - \gamma$ , that is,

$$\text{power} = 1 - \gamma = \Pr(\text{Reject } H_0 \mid H_0 \text{ is false}).$$

In simple terms, the power of a test is the probability that the study will determine that there is a treatment effect of some subject-matter importance when it truly exists. Since  $\gamma$  necessarily depends on the particular choice of value for  $\delta \neq 0$  under the alternative hypothesis, so too does the power of a test. Thus, with all other things being equal, the further the true value of  $\delta$  is from zero, the greater is the power of a test of  $H_0: \delta = 0$ .

By considering the two types of errors that can arise when conducting a statistical test, we can determine the sample size required to have some specified power to detect an effect,  $\delta \neq 0$ . For the special case of the two group comparison considered in our example, a test of  $H_0: \delta = 0$  can be based on the following  $z$ -test,

$$Z = \frac{\widehat{\delta}}{\sqrt{\text{Var}(\widehat{\delta})}} = \frac{\widehat{\mu}^{(1)} - \widehat{\mu}^{(2)}}{\sqrt{\frac{\sigma^2}{N_1} + \frac{\sigma^2}{N_2}}} = \frac{\widehat{\mu}^{(1)} - \widehat{\mu}^{(2)}}{\sqrt{\frac{\sigma^2}{N\pi(1-\pi)}}},$$

where  $\widehat{\delta} = \widehat{\mu}^{(1)} - \widehat{\mu}^{(2)}$  is the difference in sample means in the two groups, and  $\sigma^2$  is the variance of the response (assumed to be common in the two groups). A formula for the approximate total sample

size,  $N = N_1 + N_2$ , for a 2-tailed test is given by

$$(20.1) \quad g\{E(Y_{ij}|b_i)\} = \eta_{ij}^b = X'_{ij}\beta + Z'_{ij}b_i$$

where  $Z_{(1 \times \alpha/2)}$  and  $Z_{(1 \times \gamma)}$  denote the  $(1 - \alpha/2) \times 100\%$  and  $(1 - \gamma) \times 100\%$  percentiles of a standard normal distribution (e.g., the 97.5th percentile of a standard normal distribution is 1.96; or put in somewhat simpler terms, 97.5% of the area under the standard normal curve lies to the left of 1.96). When applying (20.1)  $\sigma^2$  is replaced by an estimate of the variability. Given the total projected sample size,  $N$ , the number of subjects in group 1 is  $N_1 = \pi N$  and the number of subjects in group 2 is  $N_2 = (1 - \pi)N$ ; these estimates of  $N_1$  and  $N_2$  are rounded up to the next nearest integer.

The main reason for displaying the formula given by (20.1) is to highlight its main constituents. A closer examination of this formula reveals that the determination of sample size requires that all of the following be specified:

1. significance level,  $\alpha$ ;
2. power,  $1 - \gamma$ ;
3. effect size,  $\delta$ ; and
4. common variance,  $\sigma^2$ .

Ordinarily, the first two factors do not pose a great challenge for investigators. Conventionally, the significance level of a statistical test is fixed at the mythical 0.05 level (with  $Z_{(1 \times \alpha/2)} = 1.96$  for a 2-tailed test). Similarly the lower bound on what might be considered acceptable power is usually set at approximately 80% (with  $Z_{(1 \times \gamma)} = 0.842$  for power = 0.8, or  $Z_{(1 \times \gamma)} = 1.282$  for power = 0.9). This leaves only two key ingredients for which the investigators must provide information: the minimum effect size of scientific interest and an estimate of the variability in the data. Note that the former appears in the denominator of (20.1), while the latter appears in the numerator. As a result, for any fixed value of the variability, the required sample size decreases with increasing effect size,  $\delta$ . Intuitively, fewer subjects (or less information) are needed when it is of interest to determine whether a true treatment effect is quite far from the null value. Similarly, for any fixed effect size, the required sample size decreases with decreasing variability. For example, the required sample size can be made smaller by using a more reliable measurement instrument.

## 20.3 SAMPLE SIZE FOR A LONGITUDINAL CONTINUOUS RESPONSE

We first consider the common scenario where investigators are interested in comparing two treatments, an active treatment and control, in terms of *changes* in the mean response over time. Toward the end of the section, we also consider the less common scenario where investigators are interested in the comparison of the *time-averaged* response (i.e., the *average* response over the duration of the study) rather than *changes* in the mean response. Throughout, we assume that investigators plan to randomize a total of  $N$  subjects, with  $N_1 = \pi N$  in group 1 (e.g., active treatment), and  $N_2 = (1 - \pi)N$  in group 2 (e.g., control). They plan to take  $n$  repeated measurements of the response (not necessarily equally spaced measurements).

## 20.3.1 Sample Size for Comparison of Change in Response

In this section we consider the case where, at the completion of the study, the two treatment groups are to be compared in terms of changes in the mean response over the duration of the study. For simplicity we assume that changes in the mean response can be expressed in terms of a linear trend and the treatment effect can be expressed in terms of the difference in slopes or rates of change, say  $\delta$ . Under the null hypothesis of no treatment difference, that is, no treatment  $\times$  linear trend interaction,  $H_0: \delta = 0$ . We show that sample size calculations for such a longitudinal study design can be simplified so that the standard sample size formula given by (20.1) can be used. This is achieved by considering the two-stage model for longitudinal data described in Chapter 8 (see Section 8.4). Let us assume the following two-stage formulation. At the first stage, we assume that a simple parametric curve (e.g., linear trend in time) fits the observed responses for each subject. In the second stage, these individual-specific parameters are then related to covariates that describe the different groups from which the individuals have been drawn (e.g., active treatment versus control).

**Stage 1:** In the first stage subjects are assumed to have their own unique individual-specific response trajectories. That is, in stage 1 we posit that the repeated measures on each individual follow a regression model with a linear trend in time, but with separate regression coefficients for each individual

$$Y_{ij} = \beta_{1i} + \beta_{2i} t_j + \epsilon_{ij},$$

where the errors,  $\epsilon_{ij}$ , are assumed to be independent and identically distributed, having a normal distribution with mean equal to zero and variance  $\sigma_\epsilon^2$ , that is,  $\epsilon_{ij} \sim N(0, \sigma_\epsilon^2)$ .

**Stage 2:** In the second stage we assume that the individual-specific effects,  $\beta_i = (\beta_{1i}, \beta_{2i})'$ , are random. The mean and covariance of  $\beta_{1i}$  and  $\beta_{2i}$  are the population parameters that are modeled in the second stage. Specifically, variation in  $\beta_i$  is modeled as a function of between-individual covariates, which we assume here include only the treatment group. Thus we can allow the mean of  $\beta_i$  (i.e., the mean intercept and slope) to depend on the treatment group,

$$E(\beta_{1i} | \text{Group}_i = g) = \beta_1^{(g)}, \text{ for } g = 1, 2,$$

$$E(\beta_{2i} | \text{Group}_i = g) = \beta_2^{(g)}, \text{ for } g = 1, 2,$$

where  $\text{Group}_i = 1$  if the  $i^{th}$  individual was assigned to the active treatment, and  $\text{Group}_i = 2$  otherwise. In this model,  $\beta_1^{(1)}$  is the mean slope, or constant rate of change in the mean response over time, in the active treatment group, while  $\beta_1^{(2)}$  is the mean slope in the control group. That is,  $\beta_1^{(1)} - \beta_1^{(2)}$  has interpretation in terms of a treatment group difference in the rate of change in the mean response and corresponds to the definition of  $\delta$  given earlier. The residual between-individual variation in the  $\beta_i$  that cannot be explained by treatment group is expressed as

$$\text{Cov}(\beta_i) = G = \begin{pmatrix} g_{11} & g_{12} \\ g_{21} & g_{22} \end{pmatrix},$$

where  $g_{11} = \text{Var}(\beta_{1i})$ ,  $g_{22} = \text{Var}(\beta_{2i})$ , and  $g_{12} = g_{21} = \text{Cov}(\beta_{1i}, \beta_{2i})$ .

This two-stage formulation yields a tractable form for the component of variability required for sample size and power calculations. If each subject is measured at a common set of occasions,  $t_1, \dots, t_n$ , and there are  $N_1$  and  $N_2$  subjects in the two treatment groups (for a total sample size of  $N$ ), we can derive simple expressions for sample size and power similar to the univariate setting. Letting  $\hat{\beta}_{2i}$  denote the ordinary least squares (OLS) estimate of the slope for the  $i^{th}$  subject, the variability of  $\hat{\beta}_{2i}$  is given by

$$\text{Var}(\hat{\beta}_{2i}) = \sigma_\epsilon^2 \left\{ \sum_{j=1}^n (t_j - \bar{t})^2 \right\}^{-1} + g_{22},$$

where

$$\bar{t} = \frac{1}{n} \sum_{j=1}^n t_j.$$

Thus the variability of  $\hat{\beta}_{2i}$  is composed of two components: the within-subject variance,  $\sigma_\epsilon^2 \{\sum_{j=1}^n (t_j - \bar{t})^2\}^{-1}$ , and the between-subject variance,  $g_{22} = \text{Var}(\beta_{2i})$ . To test whether the mean slopes are equal in the two treatment groups, we can construct the following  $z$ -test based on averages of the  $\hat{\beta}_{2i}$ ,

$$Z = \frac{\hat{\delta}}{\sqrt{\text{Var}(\hat{\delta})}} = \frac{\bar{\beta}_2^{(1)} - \bar{\beta}_2^{(2)}}{\sqrt{\frac{\sigma_\beta^2}{N_1} + \frac{\sigma_\beta^2}{N_1}}} = \frac{\bar{\beta}_2^{(1)} - \bar{\beta}_2^{(2)}}{\sqrt{\frac{\sigma_\beta^2}{N\pi(1-\pi)}},$$

where  $\bar{\beta}_2^{(1)}$  and  $\bar{\beta}_2^{(2)}$  are the sample averages of  $\hat{\beta}_{2i}$  in the treatment and control groups respectively,  $\sigma_\beta^2 = \text{Var}(\hat{\beta}_{2i})$ , and  $\pi$  is the proportion of subjects in Group 1.

Given estimates of  $g_{22}$ , the between-subject variability in slopes, and  $\sigma_\epsilon^2$ , the within-subject variability, the sample size can be determined from the standard formula (20.1) introduced in Section 20.2,

$$(20.2) \quad N = \frac{\{Z_{(1-\alpha/2)} + Z_{(1-\gamma)}\}^2 \sigma_\beta^2}{\pi(1-\pi)\delta^2},$$

where now

$$\sigma_\beta^2 = \text{Var}(\hat{\beta}_{2i}) = \sigma_\epsilon^2 \left\{ \sum_{j=1}^n (t_j - \bar{t})^2 \right\}^{-1} + g_{22},$$

and  $\delta$  is a treatment effect size of interest (i.e.,  $\delta$  is the treatment group difference in slopes or rates of change in the mean response). Notice that the sample size formula given by (20.2) is virtually identical to (20.1), except that  $\sigma_\beta^2$  has two components: a within-subject variance component,  $\sigma_\epsilon^2 \{\sum_{j=1}^n (t_j - \bar{t})^2\}^{-1}$ , and a between-subject variance component,  $g_{22} = \text{Var}(\beta_{2i})$ . When applying (20.2),  $\sigma_\beta^2$  is replaced by estimates of these two sources of variability. Furthermore, if the measurement occasions are equally spaced (at least approximately) throughout the duration of the study, then the expression for  $\sum_{j=1}^n (t_j - \bar{t})^2$  simplifies to  $\{\tau^2 n(n+1)\}/\{12(n-1)\}$ , where  $\tau$  denotes the duration of the study. (The latter expression can be derived by using the fact that the variance of the first  $n$  integers is  $(n+1)(n-1)/12$ .)

The sample size formula given by (20.2) can also be manipulated to determine the power of a test of  $H_0$  for a given sample size, since (20.2) implies that

$$(20.3) \quad Z_{(1-\gamma)} = \sqrt{\frac{N\pi(1-\pi)\delta^2}{\sigma_\beta^2}} - Z_{(1-\alpha/2)}.$$

Therefore the power,  $1 - \gamma$ , is given by  $\Phi\{Z_{(1-\gamma)}\}$ , where  $\Phi(\cdot)$  denotes the cumulative standard normal distribution function. That is, the value of  $Z_{(1-\gamma)}$  can be calculated from (20.3), and the power is determined by the area under the standard normal curve that lies to the left of  $Z_{(1-\gamma)}$ . For example, if  $Z_{(1-\gamma)} = 1.08$  for some given sample size, then the projected power is 0.86, corresponding to the area under the standard normal curve that lies to the left of 1.08.

For more general models (e.g., quadratic trends or a spline function of time), sample size formulas for two group comparisons of other coefficients in the model for the mean response can be derived by using the general formulation for the stage 1 model (see Section 8.4),

$$Y_i = Z_i \beta_i + \epsilon_i,$$

where the matrix  $Z_i$  specifies how an individual's responses change over time and  $\beta_i$  is a  $q \times 1$  vector of individual-specific regression coefficients. Then, for any particular trend of interest, the variance,  $\sigma_\beta^2$ , in the sample size formula is simply obtained from the appropriate diagonal element of

$$\begin{aligned}\text{Cov}(\widehat{\beta}_i) &= \sigma_\epsilon^2 (Z_i' Z_i)^{-1} + \text{Cov}(\beta_i) \\ &= \sigma_\epsilon^2 (Z_i' Z_i)^{-1} + G.\end{aligned}$$

In our previous example, with random intercepts and slopes,

$$\begin{aligned}\text{Cov}(\widehat{\beta}_i) &= \text{Cov} \left( \begin{array}{c} \widehat{\beta}_{1i} \\ \widehat{\beta}_{2i} \end{array} \right) \\ &= \sigma_\epsilon^2 \left[ \left( \begin{array}{c} 1, \dots, 1 \\ t_1, \dots, t_n \end{array} \right) \left( \begin{array}{cc} 1 & t_1 \\ \vdots & \vdots \\ 1 & t_n \end{array} \right) \right]^{-1} + \left( \begin{array}{cc} g_{11} & g_{12} \\ g_{21} & g_{22} \end{array} \right),\end{aligned}$$

and  $\text{Var}(\widehat{\beta}_{2i})$  is the lower-diagonal element (2nd row, 2nd column) of this  $2 \times 2$  matrix.

In the absence of any information about the variability of the random slopes, the simple formulas for sample size and power given by (20.2) and (20.3) cannot be used. A number of textbooks, and commercially available software for sample size and power calculations, rely on formulas that assume a longitudinal model with only randomly varying intercepts (or random subject effects). This implies a compound symmetry covariance matrix, with constant variance over time, say  $\text{Var}(Y_{ij}) = \sigma_b^2 + \sigma_\epsilon^2$  (where  $\sigma_b^2$  and  $\sigma_\epsilon^2$  denote the between-subject and within-subject variances, respectively), and constant correlation, say  $\rho$ , among pairs of repeated measurements, where

$$\rho = \text{Corr}(Y_{ij}, Y_{ik}) = \frac{\sigma_b^2}{\sigma_b^2 + \sigma_\epsilon^2}.$$

By making the strong assumption of no between-subject variability in the slopes ( $g_{22} = 0$ ), the sample size formula given by (20.2) can be used where

$$\sigma_\beta^2 = \sigma_\epsilon^2 \left\{ \sum_{j=1}^n (t_j - \bar{t})^2 \right\}^{-1} + g_{22}$$

is replaced by

$$\sigma_\beta^2 = \sigma_\epsilon^2 \left\{ \sum_{j=1}^n (t_j - \bar{t})^2 \right\}^{-1} = (1 - \rho)(\sigma_b^2 + \sigma_\epsilon^2) \left\{ \sum_{j=1}^n (t_j - \bar{t})^2 \right\}^{-1}.$$

This yields the following formula for sample size,

$$\begin{aligned}N &= \frac{\{Z_{(1-\alpha/2)} + Z_{(1-\gamma)}\}^2 \sigma_\epsilon^2}{\pi(1 - \pi) \delta^2 \sum_{j=1}^n (t_j - \bar{t})^2} \\ (20.4) \quad &= \frac{\{Z_{(1-\alpha/2)} + Z_{(1-\gamma)}\}^2 (1 - \rho) \text{Var}(Y_{ij})}{\pi(1 - \pi) \delta^2 \sum_{j=1}^n (t_j - \bar{t})^2};\end{aligned}$$

a similar formula for power can be based on (20.3) with  $\sigma_\beta^2$  replaced by the expression above. This simple formula for sample size is instructive about how sample size (and power) is impacted by the correlation among repeated measures. It is apparent from (20.4) that  $N$  decreases with increasing (positive) correlation,  $\rho$ . In general, when estimating and comparing *change* in the response over time, the sample size required decreases with increasing magnitude of the positive correlations among the responses. Although the formula given by (20.4) appears in numerous textbooks that discuss sample size (and power) calculations for longitudinal study designs, we must caution the reader that use of (20.4) can dramatically underestimate the required sample size. That is, it can be shown that the assumption of equal variances and equal correlations (compound symmetry) on which (20.4) is based always produces an underestimate of the projected sample size if a more complex random effects covariance structure (e.g., random intercepts and slopes) or an arbitrary covariance matrix actually holds.

Finally, we conclude this section by noting that the sample size formulas given by (20.2) and (20.4) are special cases of a more general formula based upon contrasts or linear summary statistics for the response vector, say

$$C_i = \sum_{j=1}^n a_j Y_{ij},$$

for a set of known weights  $a_j$  (for  $j = 1, \dots, n$ ), and where

$$\delta = E(C_i | \text{Group}_i = 1) - E(C_i | \text{Group}_i = 2).$$

For example, the ordinary least squares (OLS) estimator of the slope for the  $i^{th}$  subject, considered earlier in this section, can be expressed as

$$C_i = \sum_{j=1}^n a_j Y_{ij}, \text{ where } a_j = \frac{(t_j - \bar{t})}{\sum_{j=1}^n (t_j - \bar{t})^2};$$

other choices of weights can be used to summarize the quadratic trend over time (e.g., polynomial contrast coefficients), area under the curve (AUC), or the mean of all post-baseline measurements minus the baseline. For example, with five equally spaced measurements, a summary statistic for the quadratic time trend is given by

$$C_i = 2 \times Y_{i1} - 1 \times Y_{i2} - 2 \times Y_{i3} - 1 \times Y_{i4} + 2 \times Y_{i5},$$

with  $a_1 = 2, a_2 = -1, a_3 = -2, a_4 = -1, a_5 = 2$ . Similarly a summary statistic for the mean of all post-baseline measurements minus the baseline is given by

$$C_i = -1 \times Y_{i1} + \frac{1}{4} \times (Y_{i2} + Y_{i3} + Y_{i4} + Y_{i5}), (\text{with } a_1 = -1, a_2 = a_3 = a_4 = a_5 = \frac{1}{4}).$$

For any choice of weights used to form the summary statistic,  $C_i$ , sample size can be determined by the following general formula:

$$(20.5) \quad N = \frac{\{Z_{(1-\alpha/2)} + Z_{(1-\gamma)}\}^2 \sigma_C^2}{\pi(1-\pi)\delta^2},$$

where now

$$(20.6) \quad \sigma_C^2 = \text{Var}(C_i) = \sum_{j=1}^n \sum_{k=1}^n a_j \times a_k \times \sigma_{jk},$$

and  $\sigma_{jk} = \text{Cov}(Y_{ij}, Y_{ik})$ . This formula can be used for two-group comparisons of any summary statistic (e.g., slope, polynomial contrast, AUC minus baseline); in addition it allows for more general patterns for the covariance among the responses. In Section 20.3.3 we use this formula for the comparison of a “time-averaged” response.

## 20.3.2 Effects of Study Duration and Frequency of Observation

In this section we closely examine the simple formula for sample size given by (20.2) to reveal how sample size (and power) is impacted by:

1. the length of the study,
2. the number of repeated measures, and
3. the spacing of the repeated measures.

Consider the design of a longitudinal study with  $n$  repeated measurements. Let  $t_1 = 0$  denote the baseline measurement occasion and  $t_n = \tau$  denote the time at the final measurement occasion (i.e.,  $\tau$  denotes the duration of the study). Study investigators can reduce the required sample size (or, correspondingly, increase the power for a fixed total sample size,  $N$ ) by reducing the magnitude of  $\sigma_{\beta}^2$ . Recall that  $\sigma_{\beta}^2$  depends on both the between-subject and within-subject variability. In general, investigators have relatively little control over the natural heterogeneity of the study population, denoted in this instance by  $g_{22} = \text{Var}(\beta_{2i})$  and, more generally, by  $G = \text{Var}(\beta_i)$ . The between-subject variability can only be reduced by focusing on a more homogeneous population. However, doing so could alter the intended target of inference and reduce the generalizability of the results. In principle, the within-subject variability,  $\sigma_{\epsilon}^2$ , can be reduced by using a more reliable measurement instrument; however, this will not always be possible or practical. Therefore, to reduce the magnitude of  $\sigma_{\beta}^2$ , we must focus on ways to increase the magnitude of

$$\sum_{j=1}^n (t_j - \bar{t})^2.$$

Because  $\sum_{j=1}^n (t_j - \bar{t})^2$  is a divisor of  $\sigma_{\epsilon}^2$ , increasing its magnitude reduces the contribution of the within-subject variance to  $\sigma_{\beta}^2$ . Note that  $\sum_{j=1}^n (t_j - \bar{t})^2$  is the sum of the squared deviations of the measurement times about their mean. It is a function of the duration of the study,  $\tau$ , the number of repeated measurements,  $n$ , and the relative spacing of the repeated measurements. For a study of fixed length  $\tau$  and fixed number of repeated measures  $n$ ,  $\sum_{j=1}^n (t_j - \bar{t})^2$  is maximized when  $n/2$  measurements are taken at baseline and  $n/2$  measurements are taken at the end of the study (when  $n$  is an even number). In general, such a study design would not be desirable because it relies too heavily on the assumption that changes in the response are linear over time and precludes examination of non-linear (e.g., quadratic) trends. Also the notion of taking  $n/2$  replicate measurements at the same occasion is not feasible or practical in many settings. So, for the remainder of this discussion, we assume that the measurement occasions will be equally spaced (at least approximately) throughout the duration of the study. That is, in a study of length  $\tau$ , the  $n$  repeated measurements are to be taken at times  $t_1 = 0$ ,  $t_2 = \tau/(n-1)$ ,  $t_3 = 2\tau/(n-1)$ , ...,  $t_n = \tau$ . Recall that with equally spaced measurement it can be shown that

$$\sum_{j=1}^n (t_j - \bar{t})^2 = \frac{\tau^2 n(n+1)}{12(n-1)}.$$

Thus, for a fixed number of repeated measurements, doubling the length of the study decreases the impact of the within-subject variability by a factor of 4. Impressive as this may seem, there are a number of practical limitations that qualify this result. First, the length of a longitudinal study is usually determined by economic, logistical, and subject-matter factors that constrain the maximum length of follow-up. Second, changes in the mean response, as a function of exposure to some treatment or intervention, may be of limited duration and constrain the maximum value of  $\tau$ . As a result many study investigators are restricted to a relatively narrow range of possible values for  $\tau$ . Third, increasing the duration of a study may potentially increase the rate of attrition.

The simple formula for  $\sum_{j=1}^n (t_j - \bar{t})^2$  given above also indicates that for fixed  $\tau$  the impact of the within-subject variability decreases non-linearly with increasing  $n$ . For example, increasing the

number of repeated measurements from  $n = 2$  (a simple pre-post longitudinal design) to  $n = 4$ ,  $n = 6$ ,  $n = 8$ , and  $n = 10$ , results in a 10%, 29%, 42%, and 50% reduction in the impact of the within-subject variability. However, for a study of fixed length  $\tau$ , there may be some practical constraints on the number of repeated measurements that can be taken. Also we must caution that these results for the impact of increasing either the length of the study or the number of repeated measures rely heavily on the assumption that changes in the mean response over the duration of the study are linear in time. For example, when the assumed trend is curvilinear (e.g., quadratic trend in time), it can be shown that the number of repeated measurements,  $n$ , has an even more pronounced effect on reducing the impact of the within-subject variability.

The simple formulas for sample size and power given by (20.2) and (20.3) can be useful for making informed decisions about how best to design a longitudinal study. For any fixed values of  $\pi$ ,  $\delta$ , and  $\tau$ , the first term on the right-hand side of the formula for power (20.3),

$$\sqrt{\frac{N\pi(1-\pi)\delta^2}{\sigma_\beta^2}}$$

increases as

$$\frac{\sigma_\beta^2}{N} = \left\{ \frac{12(n-1)\sigma_\epsilon^2}{\tau^2 n(n+1)} + g_{22} \right\} / N$$

decreases. Therefore, among other factors, the power of a longitudinal study is determined by a combination of the number of subjects,  $N$ , and the number of repeated measurements,  $n$ . Increasing either the sample size or the number of repeated measurements per subject increases power. However, relative to an increase in the number of repeated measurements, increasing the sample size will generally have a far greater effect on power. The reason for this can be seen by noting that as  $n$  increases, for a fixed  $N$ , it reduces the impact of the within-subject variability ( $\sigma_\epsilon^2$ ) but leaves the between-subject variability ( $g_{22}$ ) unchanged. Therefore, even with an infinitely large number of repeated measurements, the power of the study will have an upper bound that is less than 1; how much less than 1 will be determined by the magnitude of the between-subject variability (relative to  $N$  and  $\delta$ ). In contrast, as sample size increases, for a fixed  $n$ , power increases without any bound (leading eventually to power approaching 1 for a sufficiently large  $N$ ). Put another way, an increase in the number of repeated measurements only decreases the impact of the within-subject variability, while an increase in sample size decreases the impact of both the within-subject and between-subject variability. This suggests that adding more repeated measures to a study of fixed length is most helpful only in cases where the within-subject variation is large relative to the between-subject variation. In general, adding more subjects is the most effective way to increase the power of a longitudinal study.

### 20.3.3 Sample Size for Comparison of Time-Averaged Response

So far our discussion of sample size and power has focused on study designs where we are primarily interested in comparing groups in terms of *changes* in the mean response over time. Although less common, sometimes it is of interest to design a longitudinal study where the main comparison of interest is the *average* response over the duration of the study. For example, a between-group comparison of the *time-averaged* response provides a relatively powerful analysis in a study where the treatment effect has a quick onset and numerous repeated measurements are obtained after a maximum effect has been attained. This comparison can be evaluated by calculating a single composite score (or mean) of the correlated repeated measurements for each individual, say  $C_i = \frac{1}{n} \sum_{j=1}^n Y_{ij}$ . This is an example of a linear summary statistic discussed at the end of Section 20.3.1, where  $C_i = \sum_{j=1}^n a_j Y_{ij}$  and  $a_j = \frac{1}{n}$  (for  $j = 1, \dots, n$ ). Values of  $C_i$  can then be analyzed using a standard *t*-test for independent groups. With a single summary score,  $C_i$ , sample size (and power) calculations are relatively straightforward and can rely on the general formula for the difference between means for two independent groups (see [equation \(20.5\)](#) in Section 20.3.1). Specifically, a formula for the approximate total sample size,  $N = N_1 + N_2$ , is given by

$$N = \frac{\{Z_{(1-\alpha/2)} + Z_{(1-\gamma)}\}^2 \sigma_C^2}{\pi(1-\pi)\delta^2},$$

where  $\sigma_C^2 = \text{Var}(C_i)$  denotes the variance of  $C_i$  (assumed to be common for the two groups), and  $\delta = \{E(C_i|\text{Group}_i = 1) - E(C_i|\text{Group}_i = 2)\}$  is the difference in the mean of  $C_i$  for the two groups. Recall that the variance of a sum (or average) of  $n$  repeated measurements is a function of the variances of the measurements at the  $n$  occasions and their intercorrelations. Specifically,

$$\text{Var}(C_i) = \frac{1}{n^2} \sum_{j=1}^n \sum_{k=1}^n \rho_{jk} \times \sigma_j \times \sigma_k,$$

where  $\rho_{jk} = \text{Corr}(Y_{ij}, Y_{ik})$  and  $\sigma_j^2 = \text{Var}(Y_{ij})$ . Note that this expression for  $\text{Var}(C_i)$  is a special case of [equation \(20.6\)](#) in Section 20.3.1, where  $a_j = \frac{1}{n}$  (for  $j = 1, \dots, n$ ). A simpler expression for  $\text{Var}(C_i)$  is obtained when it can be assumed that the variances are constant over time, with  $\sigma_j^2 = \sigma_k^2 = \sigma^2$ ,

$$\text{Var}(C_i) = \frac{\sigma^2}{n^2} \sum_{j=1}^n \sum_{k=1}^n \rho_{jk};$$

if in addition it is assumed that the correlations are approximately constant, with  $\text{Corr}(Y_{ij}, Y_{ik}) = \rho$ , then

$$\text{Var}(C_i) = \frac{1}{n} \{1 + (n-1)\rho\} \sigma^2.$$

The design and analysis of a study with a time-averaged response,  $C_i$ , require greater care when  $Y_{i1}$  is a “baseline” response; see Chapter 5, Sections 5.6 and 5.7, for a comparison of alternative strategies for handling the baseline response. For example, in the setting of a randomized trial, where  $Y_{i1}$  is a baseline (pre-randomization) response, the comparison of groups in terms of  $C_i$  is a less meaningful indicator of treatment effect because the expected difference between treatment groups at baseline is known to be zero. Inclusion of the baseline response in  $C_i$  would bias the comparison of the treatment groups toward the null. There are two alternative approaches to the analysis in this setting. First, the baseline response can be excluded from the construction of the time-averaged score, with

$$C_i = \frac{1}{n-1} \sum_{j=2}^n Y_{ij}$$

based only on the  $(n-1)$  post-baseline responses. Values of  $C_i$  can then be analyzed using a standard *t*-test for independent groups. Second, the analysis of  $C_i$  based on the  $(n-1)$  post-baseline responses

can be adjusted for the baseline response; that is, we can test for group differences in  $C_i$  using analysis of covariance (ANCOVA), with baseline response as a covariate. In general, the latter is the preferred method of analysis because it can substantially reduce the variability of  $C_i$ , leading to a more powerful test of group differences and a study requiring fewer subjects (and/or fewer repeated measurements). With adjustment for the baseline response,  $Y_{i1}$ , the sample size formula given above can be used by replacing  $\sigma_C^2$  with

$$\sigma_C^2 = \text{Var}(C_i) [1 - \{\text{Corr}(Y_{i1}, C_i)\}^2],$$

where  $\text{Var}(C_i)$  is defined above (albeit now for  $C_i$  based only on the  $(n - 1)$  post-baseline responses), and

$$\begin{aligned}\text{Corr}(Y_{i1}, C_i) &= \frac{\text{Cov}(Y_{i1}, C_i)}{\sqrt{\text{Var}(Y_{i1}) \text{Var}(C_i)}} \\ &= \frac{\frac{1}{n-1} \sum_{j=2}^n \rho_{1j} \times \sigma_1 \times \sigma_j}{\sqrt{\sigma_1^2 \frac{1}{(n-1)^2} \sum_{j=2}^n \sum_{k=2}^n \rho_{jk} \times \sigma_j \times \sigma_k}} \\ &= \frac{\sum_{j=2}^n \rho_{1j} \times \sigma_j}{\sqrt{\sum_{j=2}^n \sum_{k=2}^n \rho_{jk} \times \sigma_j \times \sigma_k}}.\end{aligned}$$

From the expression above for  $\sigma_C^2$  it is apparent that adjustment for baseline response reduces variability when the baseline response is correlated with  $C_i$ . This reduction in variability decreases the required sample size (or, correspondingly, increases the power for a fixed total sample size,  $N$ ). If it can be assumed that  $\sigma_j^2 = \sigma_k^2 = \sigma^2$ , and that the correlations are approximately constant, with  $\text{Corr}(Y_{ij}, Y_{ik}) = \rho$ , then  $\sigma_C^2$  in the sample size formula can be replaced by

$$\sigma_C^2 = \text{Var}(C_i) [1 - \{\text{Corr}(Y_{i1}, C_i)\}^2] = \frac{(1 - \rho)\{1 + (n - 1)\rho\} \sigma^2}{(n - 1)}.$$

Although the formula for  $\sigma_C^2$  given above is simple, we must remind the reader that the assumption of equal variances and equal correlations produces inaccurate estimates of the projected sample size if a more complex covariance matrix actually holds.

When there is no a priori reason to assume the groups have the same mean response at baseline (e.g., in an observational study), the preferred strategy for adjusting for baseline response is to compare the groups in terms of the mean of all post-baseline measurements minus the baseline (see Section 20.3.1). That is, the comparison of groups can be made in terms of  $C_i$ , where

$$C_i = \left( \frac{1}{n-1} \sum_{j=2}^n Y_{ij} \right) - Y_{i1}.$$

In concluding this section on sample size and power for a longitudinal continuous response, we note that the focus has been exclusively on the simple two-group study design. Longitudinal studies comparing three or more treatment groups are not uncommon. Although sample size and power calculations for general linear models (e.g., ANOVA with three or more groups or linear regression with a quantitative covariate) involve a level of complexity greater than that required for the simple two group setting, the key ingredients required for their computation remain the same. Using the variance formulas for linear summary statistics,  $\sigma_\beta^2$  and  $\sigma_C^2$ , outlined in earlier sections, the extensions to three or more treatment groups, or to the regression setting with a quantitative covariate, are relatively straightforward. Finally, throughout all our discussion of sample size and power, we have assumed no missing data or attrition. The impact of missing data is difficult to quantify precisely because it depends on the patterns of missingness, and in this case simple formulas no longer apply. An admittedly ad hoc, but conservative, approach for adjusting for attrition is to inflate the required sample size in each group to account for the assumed rate of attrition (or proportion of subjects who drop out before the completion of the study). That is, if the rate of attrition is assumed to be 10% in each group, then the projected total sample size should be  $N/0.9$ .

## 20.3.4 Example: Longitudinal Study with a Continuous Response

To illustrate the application of the sample size formula (20.2), let us return to the example introduced earlier. Suppose that investigators are interested in comparing two treatments, an active treatment and control, in terms of changes in the mean response over time. The investigators plan to randomize an equal number of subjects ( $N_1 = N_2$ , with  $\pi = 0.5$ ) to receive either of the two treatments. They plan to take five repeated measurements of the response, one at baseline, and the remainder at 6-month intervals until the completion of the study ( $\tau = 2$  years). The response variable is continuous and assumed to have an approximate normal distribution. At the completion of the study, the two treatment groups are to be compared in terms of changes in the mean response over the duration of the study. For simplicity we assume that changes in the mean response can be expressed in terms of a linear trend over time (in years) and the treatment effect can be expressed in terms of the difference in slopes, say  $\delta$ .

Suppose that the investigators want to detect a minimum treatment effect of  $\delta = 1.2$ , that is, a difference in the annual rates of change in the treatment and control groups of no less than 1.2. Based on historical data from similar populations, the investigators posit that the between-subject variability in the rate of change,  $\text{Var}(\beta_{2i}) \approx 2$  and the within-subject variability,  $\sigma_\epsilon^2 \approx 7$ . Finally, the investigators desire to have power of 90% when conducting a 2-sided test at the 5% significance level (i.e.,  $\gamma = 0.1$  and  $\alpha = 0.05$ ). Given these specifications,

$$\sigma_\epsilon^2 \left\{ \sum_{j=1}^n (t_j - \bar{t})^2 \right\}^{-1} = \frac{12(n-1)\sigma_\epsilon^2}{\tau^2 n(n+1)} = \frac{12 \times 4 \times 7}{4 \times 5 \times 6} = 2.8$$

and

$$\sigma_\beta^2 = \sigma_\epsilon^2 \left\{ \sum_{j=1}^n (t_j - \bar{t})^2 \right\}^{-1} + g_{22} = 2.8 + 2.0 = 4.8.$$

The projected total sample size required is

$$N = \frac{\{Z_{(1-\alpha/2)} + Z_{(1-\gamma)}\}^2 4\sigma_\beta^2}{\delta^2} = \frac{(1.96 + 1.282)^2 \times 4 \times 4.8}{1.44} = 140.1.$$

Thus, to ensure that they have power of at least 90%, the investigators will need to enroll a total of 142 subjects, randomizing an equal number (71) to each of the two treatment groups. Note that a study of the same duration ( $\tau = 2$  years) with  $n = 3$  repeated measurements, 12 months apart, would require a total of 162 subjects to achieve comparable power. Alternatively, if it were feasible to conduct the study over 3 years instead of 2 years (and have the same retention rate), with  $n = 5$  repeated measurements taken 9 months apart, it would require a total of 96 subjects to achieve power of at least 90%.

For this example it is of interest to study the: relationship of power, sample size, and the number of repeated measurements (assuming a study of the same duration  $\tau = 2$  years). [Table 20.1](#) displays the power as a function of the total sample size,  $N$ , and the number of equally spaced repeated measurements,  $n$ . [Table 20.1](#) is revealing about the trade-offs of increasing the sample size versus increasing the number of repeated measurements. For example, doubling the sample size leads to a discernibly greater increase in power than doubling of the number of repeated measurements. This can be explained by the fact that increases in the number of repeated measurements only reduce the impact of the within-subject variance component in the formula for power. Recall that  $\sigma_\beta^2$  depends on both the between-subject and within-subject variability. In contrast, increasing the sample size reduces the impact of both sources of variability.

[Table 20.1](#) Power as a function of sample size and the number of equally spaced repeated measurements in a longitudinal study of fixed duration.

Sample Size ( $N$ )	Number of Repeated Measures ( $n$ )				
	2	4	6	8	10
40	0.37	0.39	0.43	0.47	0.50
80	0.63	0.66	0.72	0.76	0.79
120	0.80	0.83	0.87	0.90	0.93
160	0.90	0.92	0.95	0.97	0.98
200	0.95	0.96	0.98	0.99	0.99

Note: Power when conducting a 2-sided test at the 5% significance level ( $\alpha = 0.05$ ) when  $\tau = 2$ ,  $\delta = 1.2$ ,  $\text{Var}(\beta_{2i}) = 2$ , and  $\sigma_e^2 = 7$ .

Next, for illustrative purposes only, we consider the projected sample size based on (20.4) under the assumption that there is no between-subject variation in the slopes ( $g_{22} = 0$ ). In that case  $\sigma_\beta^2 = 2.8$  and the projected total sample size required is

$$N = \frac{\{Z_{(1-\alpha/2)} + Z_{(1-\gamma)}\}^2 \cdot 4 \sigma_\beta^2}{\delta^2} = \frac{(1.96 + 1.282)^2 \times 4 \times 2.8}{1.44} = 81.7.$$

This suggests that the investigators will need to enroll a total of 82 subjects (in contrast to the earlier calculation of a total of 142 subjects when  $g_{22} = 2.0$ ). However, this is a substantial underestimate of the sample size required to ensure power of at least 90% if the true variability in slopes is  $g_{22} = 2.0$ . Put another way, if  $g_{22} \approx 2.0$ , then the projected sample size of 82 subjects will provide power of approximately 70% instead of 90%, resulting in a study that is under-powered to detect the effect of interest. On a somewhat technical note, it can be argued that this illustration exaggerates the potential underestimation of sample size when based on (20.4) because assuming  $g_{22} = 0$ , when in fact  $g_{22} > 0$ , is also likely to lead to a somewhat larger projected estimate of the within-subject variability, that is, a projected estimate of  $\sigma_e^2 > 7$ . While conceding this point, it does not alter the fact that use of (20.4) is problematic when  $g_{22} > 0$ . The main purpose of this illustration is to underscore our earlier warning about the use of (20.4) for sample size (and power) calculations; in general, it produces an underestimate of the projected sample size (or, correspondingly, an overestimate of the power in the sense that the projected sample size yields less power than the nominal level) when there is any between-subject variation in the slopes.

It should be apparent from (20.2) and (20.3) that sample size and power calculations are sensitive to assumptions about the covariance among the repeated measures. Because  $\sigma_\beta^2$  depends on assumptions about the magnitudes of the between-subject and within-subject variability, it is advisable to perform a sensitivity analysis to examine how sample size varies according to changes in the values of the between-subject and within-subject variances.

Finally, toward the end of the section we discussed sample size calculations for longitudinal studies where the main comparison of interest is a summary outcome. Suppose that investigators are interested in comparing two treatments, an active treatment and control, in terms of the *average* post-baseline response over time. The investigators plan to randomize an equal number of subjects ( $N_1 = N_2$ , with  $\pi = 0.5$ ) to receive either of the two treatments. They plan to take five repeated measurements of the response, one at baseline, and the remainder at 6-month intervals until the completion of the study. The response variable,  $Y_{ij}$ , is continuous and assumed to have an approximate normal distribution. At the completion of the study, the two treatment groups are to be compared in terms of the mean response over the duration of post-baseline follow-up,  $C_i = \frac{1}{4} \sum_{j=2}^5 Y_{ij}$ .

Suppose that the investigators want to detect a minimum treatment effect of  $\delta = 1.0$ , that is, an average difference between the treatment and control groups over time of no less than 1 unit. Based on historical data from similar populations, the investigators posit that the variance of the baseline response is approximately 8.0. They are less certain about the magnitude of the correlation between

pairs of repeated measures but conjecture that the correlation is likely to be in the range 0.4 to 0.8. They also conjecture than the variability of the response can be assumed to be approximately constant over the duration of the study. The investigators desire to have power of 90% when conducting a 2-sided test at the 5% significance level.

Given these specifications, it is instructive to compare the sample sizes required for (case 1) a test for group differences in  $C_i$  using a standard  $t$ -test, and (case 2) a test for group differences in  $C_i$  using ANCOVA, with baseline ( $Y_{i1}$ ) as a covariate. For this illustration we assume  $\text{Var}(Y_{ij}) = \sigma^2 = 8$ , for  $j = 1, \dots, 5$ , and  $\text{Corr}(Y_{ij}, Y_{ik}) = \rho$ , for all  $j \neq k$ , with  $\rho$  in the range 0.4 to 0.8. Using the general sample size formula for the comparison of two groups in terms of a linear summary statistic (see [equation \(20.5\)](#) in Section 20.3.1; also see Section 20.3.3), the required sample size is

$$N = \frac{\{Z_{(1-\alpha/2)} + Z_{(1-\gamma)}\}^2 \sigma_C^2}{\pi(1-\pi)\delta^2} = \frac{(1.96 + 1.282)^2 4 \sigma_C^2}{1.0} = 42.04 \sigma_C^2,$$

where

$$\sigma_C^2 = \text{Var}(C_i) = \frac{1}{n-1} \{1 + (n-2)\rho\} \sigma^2 = \frac{8}{4} (1 + 3\rho),$$

in case 1, and

$$\begin{aligned} \sigma_C^2 &= \text{Var}(C_i) [1 - \{\text{Corr}(Y_{i1}, C_i)\}^2] \\ &= \frac{(1-\rho)\{1 + (n-1)\rho\} \sigma^2}{(n-1)} = \frac{8}{4} (1-\rho)(1+4\rho), \end{aligned}$$

in case 2.

The projected total sample sizes, as a function of varying values of  $\rho$ , are presented in [Table 20.2](#). For analysis without any adjustment for baseline, the required sample size *increases* with increasing correlation. For example, with  $\rho = 0.4$ , the required sample size is 186 subjects, whereas with  $\rho = 0.8$ , the required sample size is 286 subjects, approximately 50% larger. In general, when estimating and comparing the *time-averaged* response, the sample size required increases with increasing magnitude of the positive correlations among the responses. This can be explained by the fact that the variability of the summary score is increasing in  $\rho$ . The benefits of adjustment for baseline response are apparent in [Table 20.2](#). When  $\rho = 0.5$  a total of 128 subjects are required to achieve power of at least 90%. This is almost half the number required by the analysis that fails to adjust for baseline response. The relative magnitude of the gain in power of ANCOVA over the  $t$ -test increases rapidly with  $\rho$ . When  $\rho = 0.7$  a total of 96 subjects are required; this is approximately a third of the number of subjects required by the  $t$ -test. Also, unlike the  $t$ -test, for the ANCOVA analysis the required sample size decreases with increasing correlation. This can be explained by the fact that adjusting for baseline response removes a fraction of the between-subject variability in the summary score; moreover, the fraction of variability removed from the between-subject comparison is an increasing function of  $\rho$ . In general, when the design of a longitudinal study ensures no expected difference between the groups at baseline (see Sections 5.6 and 5.7), making an adjustment for baseline response yields a substantial power advantage.

**Table 20.2** Sample size ( $N$ ), as a function of the pairwise correlation among repeated measurements, for between-group comparison of time-averaged outcome based on  $t$ -test, and ANCOVA, adjusting for the baseline response.

Method	Correlation ( $\rho$ )				
	0.4	0.5	0.6	0.7	0.8
$t$ -test	186	212	236	262	286
ANCOVA	132	128	116	96	72

*Note:* Sample size ( $N$ ) to ensure power of 0.90 when conducting a 2-sided test at the 5% significance level ( $\alpha = 0.05$ ).

## 20.4 SAMPLE SIZE FOR A LONGITUDINAL BINARY RESPONSE

Sample size determination for longitudinal studies with a binary response variable is somewhat more complicated. Complications arise from two main sources: (1) the non-linear link function (e.g., logit) usually adopted for the relationship between the mean response and covariates, and (2) the dependence of the variance on the mean. In general, simple closed-form expressions for sample size (and power), comparable to those for a continuous response, cannot be derived. Instead, precise determination of sample size and power involves more complicated procedures that, in general, require inversion of matrices. However, in this section we outline an approximate sample size (and power) formula that does yield a closed-form expression. The formula is based on the GEE approach under a “working independence” assumption (albeit with inference based on the empirical or “sandwich” variance estimator). In general, this approach can potentially yield conservative estimates of sample size, in the sense that the projected sample size may yield greater power than the nominal level; for balanced longitudinal designs, however, the formula is quite accurate.

## 20.4.1 Sample Size for Comparison of Change in Response

Similar to the derivation for a continuous response, we suppose that investigators are interested in comparing two treatments, say an active treatment and control, in terms of changes in the mean of the binary response (or probability of success) over time. The investigators plan to randomize a total of  $N$  subjects, with  $N_1 = \pi N$  in group 1, and  $N_2 = (1 - \pi)N$  in group 2. They plan to take  $n$  repeated measurements of the response at a common set of occasions,  $t_1, \dots, t_n$  (not necessarily equally spaced measurements). At the completion of the study, the two treatment groups are to be compared in terms of changes in the success probabilities over the duration of the study. For simplicity we assume that changes in the success probabilities can be expressed in terms of a linear trend in the log odds and that the treatment effect can be expressed in terms of the difference in slopes, say  $\delta$ . Specifically, we assume that

$$\text{logit}\{\Pr(Y_{ij} = 1 | \text{Group}_i = g)\} = \beta_1^{(g)} + \beta_2^{(g)}t_j, \quad g = 1, 2;$$

where  $\text{Group}_i = 1$  if the  $i^{\text{th}}$  individual was assigned to the active treatment, and  $\text{Group}_i = 2$  otherwise. In this model,  $\beta^{(1)}_2$  is the slope, or constant rate of change in the log odds of success over time, in the active treatment group, while  $\beta^{(2)}_2$  is the slope in the control group. That is,  $\beta^{(1)}_2 - \beta^{(2)}_2$  has interpretation in terms of a treatment group difference in the slopes and corresponds to the definition of  $\delta$  given earlier. Under the null hypothesis of no treatment difference, that is, no treatment  $\times$  linear trend interaction,  $H_0: \delta = \beta^{(1)}_2 - \beta^{(2)}_2 = 0$ .

Next we present expressions for sample size and power, based on GEE methods, that are similar to the corresponding formulas in the univariate setting. Letting  $\hat{\delta} = \hat{\beta}^{(1)}_2 - \hat{\beta}^{(2)}_2$  denote the GEE estimate of the difference in slopes, a test of  $H_0: \delta = 0$  can be based on the following  $z$ -test,

$$Z = \frac{\hat{\delta}}{\sqrt{\text{Var}(\hat{\delta})}} = \frac{\hat{\beta}_2^{(1)} - \hat{\beta}_2^{(2)}}{\sqrt{\text{Var}(\hat{\beta}_2^{(1)}) + \text{Var}(\hat{\beta}_2^{(2)})}},$$

where  $\text{Var}(\hat{\beta}^{(g)}_2)$  is the variance of the slope estimator for the  $g^{\text{th}}$  group. A formula for the approximate total sample size,  $N$ , is then given by

$$(20.7) \quad N = \frac{\{Z_{(1-\alpha/2)} + Z_{(1-\gamma)}\}^2 \sigma_\beta^2}{\pi(1-\pi)\delta^2},$$

where  $v_1 = N_1 \times \text{Var}(\hat{\beta}^{(1)}_2)$  and  $v_2 = N_2 \times \text{Var}(\hat{\beta}^{(2)}_2)$ . To apply this formula, we require expressions for  $v_g$ , for  $g = 1, 2$ .

Recall that  $\hat{\beta}^{(g)}_2$  is the GEE estimate of the slope in the  $g^{\text{th}}$  group. To determine  $v_g$ , we must make assumptions about the nature and magnitude of the correlations among the repeated binary responses. Let  $R$  denote an  $n \times n$  matrix of pairwise correlations among the repeated binary responses. The components of  $R$  are  $\rho_{jk} = \text{Corr}(Y_{ij}, Y_{ik})$ ;  $R$  is assumed to be the same in both groups. Then it can be shown that  $v_g$  is given by the lower-diagonal element (2nd row, 2nd column) of the following  $2 \times 2$  matrix:

$$(20.8) \quad \left[ \begin{pmatrix} \sigma_1^{(g)}, & \dots, & \sigma_n^{(g)} \\ t_1 \times \sigma_1^{(g)}, & \dots, & t_n \times \sigma_n^{(g)} \end{pmatrix} R^{-1} \begin{pmatrix} \sigma_1^{(g)} & t_1 \times \sigma_1^{(g)} \\ \vdots & \vdots \\ \sigma_n^{(g)} & t_n \times \sigma_n^{(g)} \end{pmatrix} \right]^{-1},$$

where  $\sigma_j^{(g)} = \sqrt{\mu_j^{(g)}(1 - \mu_j^{(g)})}$  denotes the standard deviation of the binary response in the  $g^{\text{th}}$  group at the  $j^{\text{th}}$  occasion,  $\mu_j^{(g)}$  denotes the probability of success in the  $g^{\text{th}}$  group at the  $j^{\text{th}}$  occasion, and  $R^{-1}$  is the inverse of the  $n \times n$  correlation matrix  $R$ . The  $2 \times 2$  matrix given by (20.8) is derived from the “model-based” variance estimator for the GEE with  $R$  as the “working correlation” assumption; see Section 13.2 in Chapter 13. Note, in general, there is no simple closed-form expression for  $v_g$  based

on (20.8) because it involves inversion of  $R$ , an  $n \times n$  matrix. Calculating the inverse of a  $2 \times 2$  matrix is relatively straightforward and there is a well-known analytic formula; however, for a matrix with three or more rows and columns, we require the use of computers to find its inverse. Therefore  $\nu_g$  can only be obtained from (20.8) with the aid of computer software with matrix algebra functions. However, in Section 20.6, using some statistical skullduggery, we describe how  $\nu_g$  can also be obtained via application of GEE to a “pseudo-dataset” that has been created to have the assumed structure for the mean response over time.

Although there is no simple expression for  $\nu_g$  based on (20.8), a closed-form expression can be obtained when  $\hat{\beta}_2^{(g)}$  is the GEE estimate of the slope under a “working independence” assumption for the correlation among the repeated binary responses. When a “working independence” assumption is made for the purpose of estimation, note that  $\nu_g$  still depends on the true correlation matrix,  $R$ . With a “working independence” GEE, the following closed-form expression for  $\nu_g$  is obtained,

$$(20.9) \quad \nu_g = \frac{\sum_{j=1}^n \sum_{k=1}^n \rho_{jk} \times \sigma_j^{(g)} \times \sigma_k^{(g)} \times (t_j - \bar{t}^{(g)}) \times (t_k - \bar{t}^{(g)})}{\{\sum_{j=1}^n (\sigma_j^{(g)})^2 \times (t_j - \bar{t}^{(g)})^2\}^2},$$

where

$$\bar{t}^{(g)} = \frac{\sum_{j=1}^n (\sigma_j^{(g)})^2 \times t_j}{\sum_{j=1}^n (\sigma_j^{(g)})^2},$$

is a weighted average of the times of measurement for the  $g^{th}$  group (with weights that depend on the variance at each occasion). Although at first glance this formula for  $\nu_g$  may appear somewhat daunting, it only involves addition, multiplication, and division of known quantities ( $t_j$ ,  $\mu_j^{(g)}$ , and  $\rho_{jk}$ ); therefore all of the required computations can be done using a pocket calculator or within a spreadsheet. When applying this formula for  $\nu_g$ , recall that

$$\mu_j^{(g)} = \Pr(Y_{ij} = 1 | \text{Group}_i = g) = \frac{\exp(\beta_1^{(g)} + \beta_2^{(g)} t_j)}{1 + \exp(\beta_1^{(g)} + \beta_2^{(g)} t_j)},$$

and  $\sigma_j^{(g)} = \sqrt{\mu_j^{(g)}(1 - \mu_j^{(g)})}$ , for  $g = 1, 2, j = 1, \dots, n$ .

We note that for many balanced longitudinal designs, the closed-form expression for  $\nu_g$  based on (20.9) is almost identical to the corresponding expression for  $\nu_g$  based on (20.8). For more complicated and unbalanced longitudinal designs,  $\nu_g$  based on (20.9) will be larger than  $\nu_g$  based on (20.8). As a result, reliance on the closed-form expression for  $\nu_g$  in these more complicated settings will yield conservative estimates of sample size, in the sense that the projected sample size may have greater power than the nominal level. However, for all practical purposes, the differences between sample size and power calculations based on (20.8) and (20.9) will be relatively small for many standard longitudinal designs. We also note that sample size calculations based on both (20.8) and (20.9) have been implemented in a publicly-available SAS macro called GEESIZE. (GEESIZE is available at <http://www.imbs-luebeck.de/imbs/de/software>, together with detailed documentation and a series of illustrative examples.) The macro, initially developed by Rochon (1998), is very versatile and can be applied to many different longitudinal designs for a range of different types of outcomes (e.g., continuous, binary, and count data). The macro can also incorporate monotone missing data patterns and allows for a flexible family of correlation structures that includes compound symmetry and first-order autoregressive correlation.

## 20.4.2 Sample Size for Comparison of Time-Averaged Response

So far our discussion of sample size and power has focused on study designs where we are primarily interested in comparing groups in terms of *changes* in the log odds over time. Although less common, sometimes it is of interest to design a longitudinal study where the main comparison of interest is the *average* response probability over the duration of the study. This between-group contrast can be represented in the following marginal logistic regression model,

$$\text{logit}\{\Pr(Y_{ij} = 1 | \text{Group}_i = g)\} = \beta_1^{(g)}, \quad g = 1, 2;$$

where  $\text{Group}_i = 1$  if the  $i^{\text{th}}$  individual was assigned to the active treatment, and  $\text{Group}_i = 2$  otherwise. In this model,  $\beta_1^{(g)}$ , is the log odds of success (assumed constant over the duration of the study) in the  $g^{\text{th}}$  group. The probability of success in the  $g^{\text{th}}$  group is

$$\mu^{(g)} = \frac{\exp(\beta_1^{(g)})}{1 + \exp(\beta_1^{(g)})}, \quad g = 1, 2,$$

and does not depend on measurement occasions. Here  $\beta_1^{(1)} - \beta_1^{(2)}$  has interpretation in terms of the log odds ratio and corresponds to the definition of  $\delta$  given earlier. Under the null hypothesis of no treatment difference,  $H_0: \delta = \beta_1^{(1)} - \beta_1^{(2)} = 0$ . Sample size (and power) calculation is far more straightforward in this setting because the variance of the response,  $\mu^{(g)}(1 - \mu^{(g)})$ , no longer depends on measurement occasions. Specifically, we can rely on the following sample size formula:

$$N = \frac{\{Z_{(1-\alpha/2)} + Z_{(1-\gamma)}\}^2 \sigma_\beta^2}{\delta^2},$$

where

$$\delta = \text{logit}(\mu^{(1)}) - \text{logit}(\mu^{(2)}) = \log\left(\frac{\mu^{(1)}}{1 - \mu^{(1)}} / \frac{\mu^{(2)}}{1 - \mu^{(2)}}\right),$$

is the log odds ratio comparing the two groups,

$$\sigma_\beta^2 = \left(\frac{1}{n^2} \sum_{j=1}^n \sum_{k=1}^n \rho_{jk}\right) \left\{ \frac{1}{\mu^{(1)}(1 - \mu^{(1)})\pi} + \frac{1}{\mu^{(2)}(1 - \mu^{(2)})(1 - \pi)} \right\},$$

and  $\rho_{jk} = \text{Corr}(Y_{ij}, Y_{ik})$ . When it can be assumed that the correlations are approximately constant, with  $\text{Corr}(Y_{ij}, Y_{ik}) = \rho$ , for all  $j \neq k$ , then a simpler expression for  $\sigma_\beta^2$  is obtained:

$$\sigma_\beta^2 = \frac{1}{n} \{1 + (n - 1)\rho\} \left\{ \frac{1}{\mu^{(1)}(1 - \mu^{(1)})\pi} + \frac{1}{\mu^{(2)}(1 - \mu^{(2)})(1 - \pi)} \right\}.$$

However, we remind the reader that assuming equal correlations produces inaccurate estimates of the projected sample size if a more complex correlation structure actually holds.

## 20.4.3 Example: Longitudinal Study with a Binary Response

To illustrate the application of the sample size formula with binary responses, let us return to the example introduced earlier. Suppose that investigators are interested in comparing two treatments, an active treatment and control, in terms of *changes* in the probability of a binary response over the duration of the study. The investigators plan to randomize an equal number of subjects ( $N_1 = N_2$ ) to receive either of the two treatments. They plan to take five repeated measurements of the response, one at baseline, and the remainder at 6-month intervals until the completion of the study ( $\tau = 2$  years), that is,  $(t_1, t_2, t_3, t_4, t_5) = (0, 0.5, 1, 1.5, 2)$ . We assume that changes in the log odds of response can be expressed (approximately) in terms of a linear trend; the treatment effect is the difference in slopes, say  $\delta$ .

The investigators assume that the baseline probability of response is approximately 0.3 for both treatment groups. At the end of two years of follow-up, they assume that the probability of response will be relatively unchanged in the control group (0.3) but will be 0.15 in the active treatment group. Based on historical data from similar populations the investigators posit that the correlation among pairs of responses is approximately 0.5. Finally, the investigators desire to have power of 90% when conducting a 2-sided test at the 5% significance level (i.e.,  $\gamma = 0.1$  and  $\alpha = 0.05$ ).

To calculate the sample size, we need to determine  $\delta$  and the response probabilities at each occasion in the two groups. Because the two groups are assumed to have the same baseline probability of response, they have common intercepts (when  $t_1 = 0$ ),

$$\beta_1^{(1)} = \beta_1^{(2)} = \text{logit}(\mu_1^{(1)}) = \log(0.30/0.70) = -0.8473.$$

The slope for time (on the log odds scale) in group 1 is

$$\begin{aligned}\beta_2^{(1)} &= \frac{\log\{\mu_5^{(1)}/(1 - \mu_5^{(1)})\} - \log\{\mu_1^{(1)}/(1 - \mu_1^{(1)})\}}{\tau} \\ &= \frac{\log(0.15/0.85) - \log(0.30/0.70)}{2} = -0.4437.\end{aligned}$$

The corresponding slope for time in group 2 is

$$\begin{aligned}\beta_2^{(2)} &= \frac{\log\{\mu_5^{(2)}/(1 - \mu_5^{(2)})\} - \log\{\mu_1^{(2)}/(1 - \mu_1^{(2)})\}}{\tau} \\ &= \frac{\log(0.30/0.70) - \log(0.30/0.70)}{2} = 0.\end{aligned}$$

Therefore  $\delta = \beta_2^{(1)} - \beta_2^{(2)} = -0.4437$ . The intercepts and slopes for the two groups can also be used to derive the response probabilities at each occasion,

$$\mu_j^{(g)} = \frac{\exp(\beta_1^{(g)} + \beta_2^{(g)} t_j)}{1 + \exp(\beta_1^{(g)} + \beta_2^{(g)} t_j)}, \quad g = 1, 2, \quad j = 1, \dots, 5.$$

For those in group 1,

$$(\mu_1^{(1)}, \mu_2^{(1)}, \mu_3^{(1)}, \mu_4^{(1)}, \mu_5^{(1)}) = (0.300, 0.256, 0.216, 0.181, 0.150),$$

while for those in group 2,

$$(\mu_1^{(2)}, \mu_2^{(2)}, \mu_3^{(2)}, \mu_4^{(2)}, \mu_5^{(2)}) = (0.30, 0.30, 0.30, 0.30, 0.30).$$

Given these specifications for  $\mu^{(g)}$ ,  $t_j$ , and  $\rho_{jk}$ , we can calculate  $v_1$  and  $v_2$  from the formula,

$$v_g = \frac{\sum_{j=1}^n \sum_{k=1}^n \rho_{jk} \times \sigma_j^{(g)} \times \sigma_k^{(g)} \times (t_j - \bar{t}^{(g)}) \times (t_k - \bar{t}^{(g)})}{\{\sum_{j=1}^n (\sigma_j^{(g)})^2 \times (t_j - \bar{t}^{(g)})^2\}^2},$$

where

$$\bar{t}^{(g)} = \frac{\sum_{j=1}^n (\sigma_j^{(g)})^2 \times t_j}{\sum_{j=1}^n (\sigma_j^{(g)})^2}; \quad \sigma_j^{(g)} = \sqrt{\mu_j^{(g)}(1 - \mu_j^{(g)})}.$$

Here  $\rho_{jk} = \rho = 0.5$  for  $j \neq k$  and  $\rho_{jk} = 1$  for  $j = k$ . This yields  $t^{(1)} = 0.8773$ ,  $t^{(2)} = 1.0$ ,  $v_1 = 1.2676$ ,

and  $\nu_2 = 0.9524$ . Therefore the projected total sample size is

$$\begin{aligned} N &= \frac{\{Z_{(1-\alpha/2)} + Z_{(1-\gamma)}\}^2 \{\nu_1/\pi + \nu_2/(1-\pi)\}}{\delta^2}, \\ &= \frac{(1.96 + 1.282)^2 (1.2676/0.5 + 0.9524/0.5)}{(-0.4437)^2} = 237.05. \end{aligned}$$

Thus to ensure that they have power of at least 90%, the investigators will need to enroll a total of 238 subjects, randomizing an equal number (119) to each of the two treatment groups.

Finally, note that a correlation of  $\rho = 0.5$  is relatively high for repeated binary responses. For example, given marginal probabilities ranging from 0.15 to 0.30, the maximum possible correlation between any pair of binary responses is constrained to be less than 0.64 (recall that with binary data, the correlations are restricted to ranges that are determined by the marginal probabilities of success; see Section 12.2). It is instructive to re-calculate sample size for smaller values of  $\rho$ . For example, with  $\rho = 0.3$ , the investigators will need to enroll a total of 328 subjects to ensure power of at least 90%. Thus, it is apparent that sample size requirements are sensitive to the magnitude of the assumed correlation among the repeated binary responses.

## 20.5 SUMMARY

In this chapter we have shown that sample size determination for longitudinal studies can often be simplified so that well-established formulas for the univariate case can be applied. For the investigators the main challenges are in the specification of the minimum effect size of subject-matter interest and in providing a realistic estimate of the anticipated variability in the measurements. The choice of an appropriate effect size must be made on purely subject-matter grounds. If the investigators expect a large effect, then it is likely to be detected with a relatively small sample size. In contrast, detection of small effects requires somewhat larger sample sizes. In planning a study, investigators need to keep their optimism in check, since gross overestimation of the effect size will result in too few subjects and insufficient power to detect somewhat smaller, but nonetheless scientifically important, effects.

Perhaps the greatest challenge facing investigators is to provide a realistic estimate of the variability in the data. This will either require the provision of estimates of both between-subject and within-subject variability or, alternatively, an estimate of the covariance among the repeated measurements. Since scientific studies are rarely conducted in a vacuum, investigators can usually obtain some estimates of the variability based on historical data from related studies with similar populations. Alternatively, in the complete absence of any relevant historical data, it may be prudent to conduct a small pilot study. If there is much uncertainty regarding the anticipated variability in the data, a simple sensitivity analysis, examining the projected sample sizes across a range of plausible values for the variability, should be conducted.

Finally, as mentioned in Chapters 17 and 18, missing data are the rule, not the exception, in longitudinal studies. Therefore it is important to make some adjustment for the potential loss of information due to missing data when planning a longitudinal study, for example, by using relatively conservative estimates of sample size. In general, a consideration of the anticipated fraction of missing data (or the proportion of subjects who drop out before the completion of the study), say  $f$ , suggests that the sample size should be inflated by a factor of  $\frac{1}{1-f}$ . That is, if 15% of the observations are expected to be missing, then investigators should plan on increasing the sample size of the study by a factor of 1.18 (or  $\frac{1}{1-0.15}$ ). Although this adjustment is crude, ignoring both the location of the missing observations and the correlation among repeated measurements, it will probably be adequate for most practical purposes. Failure to make any adjustment for missing data will result in an underestimation of the number of subjects required to attain the desired level of power.

# 20.6 COMPUTING: SAMPLE SIZE CALCULATION FOR A LONGITUDINAL BINARY RESPONSE USING PSEUDO-DATA

Recall from Section 20.4 that to apply the simple formula for sample size (and power) for a longitudinal binary response,

$$N = \frac{\{Z_{(1-\alpha/2)} + Z_{(1-\gamma)}\}^2 \{\nu_1/\pi + \nu_2/(1-\pi)\}}{\delta^2},$$

we require an estimate of  $\nu_g$ , for  $g = 1, 2$ . In general, there is no closed-form expression for  $\nu_g$  based on (20.8); instead,  $\nu_g$  can only be obtained from (20.8) with the aid of computer software with matrix algebra functions. In this section we demonstrate how  $\nu_g$  can also be obtained via application of GEE (with “working correlation” matrix fixed at the true correlations) to a “pseudo–dataset” that has been created to have the assumed structure for the mean response over time.

Before discussing the construction of the “pseudo–dataset,” and to provide some motivation for the method to be described, we note that the required term in the numerator of the sample size formula,  $\{\nu_1/\pi + \nu_2/(1-\pi)\}$ , can be re-expressed as

$$\begin{aligned} \{\nu_1/\pi + \nu_2/(1-\pi)\} &= \frac{N_1 \text{Var}(\hat{\beta}_2^{(1)})}{\pi} + \frac{N_2 \text{Var}(\hat{\beta}_2^{(2)})}{(1-\pi)} \\ &= N \text{Var}(\hat{\beta}_2^{(1)}) + N \text{Var}(\hat{\beta}_2^{(2)}) \\ &= N \text{Var}(\hat{\delta}). \end{aligned}$$

Because  $\text{Var}(\hat{\delta})$  is the variance of the estimator of  $\delta$  based on  $N$  subjects,  $\{\nu_1/\pi + \nu_2/(1-\pi)\}$  can be thought of as the variance of the estimator of  $\delta$  if based on a single representative subject instead of a sample of  $N$  subjects. To be representative of both groups, this single subject belongs to group 1 with weight  $\pi$  and belongs to group 2 with weight  $(1-\pi)$ . By creating a “pseudo–dataset” that contains repeated measures on a single representative subject,  $\{\nu_1/\pi + \nu_2/(1-\pi)\}$  can then be obtained by squaring the reported standard error for the estimate of  $\delta$  obtained from the analysis of the “pseudo–dataset”. Specifically, for each group we create  $n$  “pseudo-observations” for the  $n$  repeated measures, where the response variable is set equal to the mean at the  $n$  occasions. The  $n$  “pseudo-observations” for group 1 receive weight of  $\pi$ , whereas the  $n$  “pseudo-observations” for group 2 receive weight of  $(1-\pi)$ . In the “pseudo–dataset” we also include the covariates required for the planned analysis, such as an indicator of group and a variable for the time of measurement. By including a weight variable, set equal to  $\pi$  for the  $n$  observations in group 1 and to  $(1-\pi)$  for the  $n$  observations in group 2, we can allow for possible unequal allocation of subjects to the two groups.

The “pseudo–dataset” for the example of a longitudinal study with a binary response from Section 20.4 is displayed in [Table 20.3](#). Recall that in this example there are two treatment groups with equal allocation of subjects to each group ( $\pi = 0.5$ ). Five repeated measurements of the binary response are planned, one at baseline, and the remainder at 6-month intervals until the completion of the study, i.e.,  $(t_1, t_2, t_3, t_4, t_5) = (0, 0.5, 1, 1.5, 2)$ . The baseline probability of response is assumed to be 0.3 for both treatment groups, yielding common logistic regression intercepts  $\beta^{(1)}_1 = \beta^{(2)}_1 = -0.8473$ . At the end of two years of follow-up, it is assumed that the probability of response will be unchanged in the control group (0.3), but will be 0.15 in the active treatment group. This yields a slope for time in group 1 of  $\beta^{(1)}_2 = -0.4437$ , whereas the corresponding slope in group 2 is  $\beta^{(2)}_2 = 0$ ; thus  $\delta = \beta^{(1)}_2 - \beta^{(2)}_2 = -0.4437$ . The intercepts and slopes for the two groups can be used to derive the response probabilities at the five occasions. For those in group 1,

$$(\mu_1^{(1)}, \mu_2^{(1)}, \mu_3^{(1)}, \mu_4^{(1)}, \mu_5^{(1)}) = (0.3000, 0.2556, 0.2157, 0.1805, 0.1500),$$

while for those in group 2,

$$(\mu_1^{(2)}, \mu_2^{(2)}, \mu_3^{(2)}, \mu_4^{(2)}, \mu_5^{(2)}) = (0.30, 0.30, 0.30, 0.30, 0.30).$$

**Table 20.3** Illustrative commands in SAS for inputting variables from a “pseudo–dataset”.

```
DATA=pseudo;  
INPUT id wt group time y;  
one=1;  
DATALINES;  
1 0.5 1 0 0.30  
1 0.5 1 0.5 0.2555697756  
1 0.5 1 1 0.2156921486  
1 0.5 1 1.5 0.1805278643  
1 0.5 1 2 0.15  
2 0.5 2 0 0.30  
2 0.5 2 0.5 0.30  
2 0.5 2 1 0.30  
2 0.5 2 1.5 0.30  
2 0.5 2 2 0.30  
;
```

The means at each occasion in the two groups are the “pseudo-responses” in [Table 20.3](#). Associated with each “response” are the covariates and the weight variable (set equal to  $\pi$  for observations in group 1 and to  $(1 - \pi)$  for observations in group 2).

By conducting a GEE analysis of the “pseudo–dataset,” under a fixed structure for the correlation among the repeated measures, we obtain both an estimate of  $\delta$  and its standard error. The square of the standard error provides an estimate of  $\{\nu_1/\pi + \nu_2/(1 - \pi)\}$  that can then be plugged into the sample size (and power) formula. Illustrative SAS commands for using PROC GENMOD to fit a marginal logistic regression model to the “pseudo-observations” are presented in [Table 20.4](#). Note that because the response variable is a proportion rather than a binary response, we must use the events/trials syntax (albeit setting the number of trials equal to 1); see Section 13.6 for more details on the command syntax for PROC GENMOD. PROC GENMOD in SAS allows the “working correlation” matrix, R, to be fixed; here, it is assumed that there is constant correlation among the five binary responses, with  $\rho = 0.5$  (this assumption can easily be relaxed and a more general correlation structure can be specified). Finally, because the “working correlation” matrix is assumed to be correctly specified, we request that standard errors for the estimated regression parameters are obtained using the “model-based” (MODELSE) variance estimator. We caution that standard errors based on the “sandwich” estimator, the default for many software packages, cannot be used for estimating  $\{\nu_1/\pi + \nu_2/(1 - \pi)\}$ .

**Table 20.4** Illustrative commands for fitting a marginal logistic regression model, with fixed within-subject correlations, to the “pseudo-observations” using PROC GENMOD in SAS.

```
PROC GENMOD DATA=pseudo DESCENDING;  
CLASS id group;  
MODEL y/one=group time group*time / DIST=BIN LINK=LOGIT SCALE=1;  
WEIGHT wt;  
REPEATED SUBJECT=id / MODELSE  
TYPE=FIXED(1 .5 .5 .5  
          .5 1 .5 .5  
          .5 .5 1 .5  
          .5 .5 .5 1);
```

The results of fitting a marginal logistic regression model, using GEE with a fixed correlation matrix, to the “pseudo-observations” are presented in [Table 20.5](#). The estimated intercept of  $-0.8473$  corresponds to the common baseline log odds of success in the two groups (and hence the estimated group effect is zero). The estimated effect of time, the slope in group 2, is zero because, by construction, the log odds are constant over time in group 2. The estimated group  $\times$  time interaction effect,  $-0.4437$ , corresponds to the difference in slopes in the two groups; this is an estimate of  $\delta$ . The corresponding standard error of the estimate of  $\delta$ ,  $2.1071$ , provides an estimate of  $\sqrt{\nu_1/\pi + \nu_2/(1-\pi)}$ . That is,  $\{\nu_1/\pi + \nu_2/(1-\pi)\} = (2.1071)^2 = 4.4399$ . We can then plug this value into the sample size formula to obtain

$$\begin{aligned} N &= \frac{\{Z_{(1-\alpha/2)} + Z_{(1-\gamma)}\}^2 \{\nu_1/\pi + \nu_2/(1-\pi)\}}{\delta^2}, \\ &= \frac{(1.96 + 1.282)^2 (4.4399)}{(-0.4437)^2} = 237.04. \end{aligned}$$

**Table 20.5** Estimated regression coefficients and model-based standard errors from the analysis of the “pseudo-observations” in [Table 20.3](#).

Variable	Estimate	SE	Z
Intercept	-0.8473	2.7603	-0.31
I(Group = 1)	0.0000	3.9158	0.00
$t_j$	0.0000	1.3801	0.00
I(Group = 1) $\times$ $t_j$	-0.4437	2.1071	-0.21

Thus, to ensure that they have power of at least 90%, the investigators will need to enroll a total of 238 subjects, randomizing an equal number (119) to each of the two treatment groups. Recall that this is the same projected number of subjects obtained using the closed-form expression based on [\(20.9\)](#) instead of [\(20.8\)](#).

Finally, using the pseudo-data it is relatively straightforward to assess how sample size (and power) are sensitive to assumptions about the strength (or nature) of the correlation among the repeated measures. To do so only requires fixing the correlation matrix to a set of alternative values. For example, if the correlation is assumed to be constant but with  $\rho = 0.3$ , this yields a model-based standard error of the estimate of  $\delta$  of 2.4783. When this value is squared and plugged into the sample size formula,

we obtain

$$N = \frac{(1.96 + 1.282)^2 (6.1420)}{(-0.4437)^2} = 327.91.$$

Thus, to ensure that they have power of at least 90%, the investigators would need to enroll a total of 328 subjects if the correlation is 0.3 instead of 0.5.

## 20.7 FURTHER READING

Schlesselman (1973a, b), in two companion papers on the design of longitudinal studies, discusses sample size calculations and issues concerning study duration and the frequency of measurement; also see Raudenbush and Liu (2001). Muller et al. (1992) provide a comprehensive approach to the calculation of statistical power for longitudinal studies with a continuous outcome. Snijders and Bosker (1993) present sample size formulas for mixed effects models for longitudinal data. Overall and Doyle (1994) provide sample size formulas based on simple composites or contrasts among the repeated measurements. Hedeker et al. (1999) discuss sample size estimation for longitudinal study designs with a continuous outcome and allow for attrition and a variety of covariance structures for the repeated measurements. Basagana and Spiegelman (2010) discuss power and sample size calculations for longitudinal studies estimating the effect of a time-varying covariate.

The closed-form expression for sample size estimation for longitudinal study designs with a binary outcome presented in Section 20.4 is derived in Jung and Ahn (2005); they also derive a more general formula incorporating missing data patterns.

## Bibliographic Notes

Lipsitz and Fitzmaurice (1994) describe a method, based on generalized least squares, for calculating sample size for longitudinal studies with binary responses. Pan (2001) derives explicit formulas for sample size and power calculations for two-group studies with correlated binary responses. A very general method for computing sample size and statistical power for longitudinal studies, based on the generalized estimating equations approach, has been developed by Liu and Liang (1997); this method does not, in general, yield closed-form expressions for sample size and power, but the method can be implemented numerically.

# *Chapter 21*

## *Repeated Measures and Related Designs*

### **21.1 INTRODUCTION**

In this chapter we discuss the application of methods for longitudinal data to closely related study designs. In these settings individuals have multiple commensurate measurements made under different circumstances and possibly also at different times. However, the major interest of the analysis is not in changes in the response over time, but in the effect of different circumstances of measurement and/or the effects of covariates on the responses.

The first design that we will consider is the classical repeated measures design. In this setting each subject is measured under a fixed number of different conditions, often corresponding to different treatments. Interest centers on comparing the effects of the different experimental conditions on the outcome. Similar to time in a longitudinal study, experimental condition is a within-subject factor and the conditions are compared using within-subject contrasts. Such designs are popular because subject-to-subject differences in outcome are accounted for in the design. Since each subject acts as his or her own control, comparisons of the outcome under different experimental conditions are estimated free of any between-subject variation in the outcome. As a result the design is potentially very efficient relative to comparing the different experimental conditions on different groups of subjects.

The second design is one that produces what we refer to as “multiple source” data. In this setting the primary outcome of interest is measured by more than one instrument or rater. This frequently happens when the outcome is difficult to measure and is thus determined under multiple different circumstances. In Section 21.3 we describe an example in the context of measuring psychopathology in children. In this context there may be some interest in comparing the different measures, but ordinarily the main focus of the analysis is the effects of subject-specific covariates on the outcome. Hence, unlike a typical longitudinal study and also unlike a classical repeated measures design, the main interest centers on the effects of subject-specific variables on response, and possibly also their interaction with the multiple sources. In many settings the fact that there are multiple sources could be regarded as a “nuisance” feature of the study design.

The distinction we make between repeated measures and multiple source data is based on what is of primary interest in the analysis. Both share the same analytical methods for regression models with correlated data. Sometimes, however, this distinction is blurred. For example, Hernández et al. (2000) report on a validity study of a new questionnaire designed to measure physical activity and inactivity in school children in Mexico. A self-reported questionnaire was administered to both the mother and the child on two different occasions; in addition a 24-hour recall (considered the best measure, but only limited to a single day) was administered on each occasion. The average of the two 24-hour recalls was considered the “gold standard” for the analysis. Here the objective of the analysis was two-fold: to compare mean responses of the two child and two mother assessments to the average of the two 24-hour recalls, and to look at the correlations between these measures. This is clearly a multiple source data set, since the mother and child questionnaires were both intended to measure the child’s average activity levels over the period. However, the primary focus of the analysis was comparing the multiple reports to each other and to the “gold standard” and examining the correlations; here the subject-specific variables, age, gender, and socioeconomic level, were used to adjust the means and the correlations for between-subject differences. Thus, in this example, the analytic goal of comparing means is closer to a repeated measures analysis than a multiple source analysis, although comparing correlations is more typical in the multiple source analysis.

We will first describe the main features of these two designs in greater detail, and then provide some examples to illustrate the application of regression methods for correlated data to repeated measures and multiple source data.

## 21.2 REPEATED MEASURES DESIGNS

Repeated measures designs are frequently encountered in applications. In the experimental context the repeated measures design is also sometimes called a randomized block design. In the simplest setting subjects each receive  $n$  treatments or experimental conditions, and the outcome is recorded for each condition. Thus each subject has a vector of  $n$  measurements,  $Y_i = (Y_{i1}, \dots, Y_{in})'$ . The treatments may be given sequentially in a randomly assigned order, but in some settings they can be given simultaneously. An example of the latter is a classic study of topical treatments for leprosy, where each patient was given four different treatments simultaneously at four different locations on their body. After a number of days, the skin lesions were recorded at each of the four locations. Letting  $Y_{ij}$  denote the response to the  $j^{th}$  treatment ( $j = 1, \dots, 4$ ), the primary interest is in comparing the four treatments. However, the analysis must account for the fact that the observations on different treatments,  $Y_i = (Y_{i1}, \dots, Y_{i4})'$ , are correlated.

As mentioned earlier, repeated measures designs are popular because subject differences in outcome are accounted for in the design. By removing between-subject variation in the outcome from the comparisons of different experimental conditions, the repeated measures design can be very efficient relative to comparing the different experimental conditions on different groups of subjects. To illustrate this point, consider a simple repeated measures design with  $N$  subjects receiving  $n = 2$  treatment conditions (producing a total of  $2N$  observations). Let  $Y_{i1}$  denote the response for the  $i^{th}$  subject under the first condition and  $Y_{i2}$  denote the response under the second condition. A natural estimate of the effect of treatment on the mean response is

$$\hat{\delta} = \hat{\mu}_1 - \hat{\mu}_2,$$

where

$$\hat{\mu}_j = \frac{1}{N} \sum_{i=1}^N Y_{ij}.$$

The variability of this estimator of the effect of treatment is given by

$$\text{Var}(\hat{\delta}) = \text{Var} \left\{ \frac{1}{N} \sum_{i=1}^N (Y_{i1} - Y_{i2}) \right\} = \frac{1}{N} (\sigma_1^2 + \sigma_2^2 - 2\sigma_{12}),$$

where  $\sigma_j^2 = \text{Var}(Y_{ij})$ ,  $\sigma_{12} = \text{Cov}(Y_{i1}, Y_{i2}) = \rho\sigma_1\sigma_2$ , and  $\rho = \text{Corr}(Y_{i1}, Y_{i2})$ . To simplify, we assume that treatment may have an impact on the mean response, but not on the variance; this assumption can be justified when treatments are allocated within subjects randomly by time. Then,  $\sigma_1^2 = \sigma_2^2 = \sigma^2$ , and the variability of this estimator simplifies to

$$\text{Var}(\hat{\delta}) = \frac{2}{N} \{ \sigma^2(1 - \rho) \}.$$

Next consider a two-group study design, where the treatment conditions are randomized to  $2N$  subjects drawn from the population, with  $N$  subjects allocated to one condition, and the remaining  $N$  subjects allocated to the other condition (producing a total of  $2N$  observations). The effect of treatment on the mean response can be estimated by comparing the mean response in the two groups of subjects. The natural estimate of treatment effect is the same as before:

$$\hat{\gamma} = \hat{\mu}_1 - \hat{\mu}_2,$$

where

$$\hat{\mu}_j = \frac{1}{N} \sum_{i=1}^N Y_{ij}$$

is the sample mean response in the group of subjects that were randomized to the  $j^{th}$  treatment condition. The variability of this estimator of the effect of treatment is given by

$$\text{Var}(\hat{\gamma}) = \frac{1}{N} (\sigma_1^2 + \sigma_2^2),$$

where  $\sigma_j^2 = \text{Var}(Y_{ij})$ . Again, if we assume that treatment may have an impact on the mean response, but not on the variance, then  $\sigma_1^2 = \sigma_2^2 = \sigma^2$ , and

$$\text{Var}(\hat{\gamma}) = \frac{2\sigma^2}{N}.$$

The potential gain in efficiency of the repeated measures design is quantified by taking the ratio of the variances of the two estimators of the effect of treatment:

$$\frac{\text{Var}(\hat{\delta})}{\text{Var}(\hat{\gamma})} = (1 - \rho).$$

Provided that the correlation among the repeated measures is positive, the repeated measures design provides a more precise estimate of the effect of treatment. For example, when  $\rho = 0.75$ , it is four times more efficient than the two-group design or, put another way, the repeated measures design only requires a quarter of the number of observations to attain the same level of precision as the two-group design.

There are many variations on the repeated measures design. If, for example, treatments are given sequentially in a random order, it is common to use a restricted randomization scheme to balance the treatments over time. This would allow equal numbers of each treatment to be assigned at each time point, for example, with two treatments denoted A and B, half the subjects would be assigned to AB, and half to BA. Thus systematic time or sequence differences will be eliminated, and the design is potentially more efficient. In this case the number of subjects is a multiple of  $n$ , the number of treatment conditions. A special case of these designs is the crossover design, where each subject receives each treatment in a random order.

A closely related design is the split-plot design. Here there are two treatment factors, a so-called main plot factor and a sub-plot factor. When subjects are the main plots, the main plot factor is a between-subject factor and each subject receives only one level of this factor. The sub-plot factor is a within-subject factor and each subject receives all levels of the within-subject factor. Thus, in the split-plot design, one factor is randomized within subjects, just as in the repeated measures design, and the other factor is randomized between subjects. For example, in a study comparing the effects of different antibiotics (drugs) and topical gels (gels) for the treatment of eye infections, subjects (main plots) are randomized to receive one of three different oral antibiotics. The main plot factor is drug (antibiotics). In addition each subject is required to apply two different topical gels directly to the left and right eye; the two different gels are randomized to the left and right eyes. Here, the sub-plot factor is gels. In the split-plot design, both the main effects and the interaction of the two factors are usually of interest, although this design provides more precise information about the within-subject factor than the between-subject factor.

Much of the literature on repeated measures is in the context of designed experiments where treatments are allocated within subjects randomly by time or location. This has several ramifications for our discussion. First, because the factors in a randomized design are typically categorical, the analysis of these designs is ordinarily presented in the context of analysis of variance (ANOVA) rather than a general regression model with correlated data. However, the analysis of variance model can be viewed as a special case of the general linear regression model for correlated data presented in Part II. Hence the regression models for correlated data apply quite straightforwardly to the classical repeated measures design.

Second, with randomization, arguments can be made that allow one to simplify the analysis of repeated measures data, especially with balanced and complete designs. There are two main approaches to the analysis of repeated measures data (see Section 3.6), repeated measures analysis by ANOVA and repeated measures analysis by multivariate ANOVA (MANOVA). With the former, one assumes that any contrast between any two repeated measures on the same subject, say  $Y_{ij} - Y_{ik}$ , has the same variance; that is,  $\text{Var}(Y_{ij} - Y_{ik})$  is constant for all choices of  $i$  and  $j \neq k$ . For example, if the covariance matrix of the vector of repeated observations,  $Y_i$ , takes on a compound symmetry form, then the requirement for the repeated measures analysis by ANOVA is satisfied. In contrast, the model for repeated measures analysis by MANOVA allows the vector of repeated observations to have an arbitrary covariance structure, but the standard analysis is usually limited to balanced and complete designs. In addition the analysis of repeated measures by ANOVA can be considerably more powerful than the analysis by MANOVA (when the assumptions for the former hold).

If the within-subject factor is randomly allocated to subjects, then randomization arguments can be made to show that the constant variance condition on the contrast does hold. More generally, one can show that  $\text{Var}(Y_{ij})$  is constant for all  $i$  and  $j$  and  $\text{Cov}(Y_{ij}, Y_{ik})$  is constant for all  $i$  and  $j \neq k$ . The general approach to a randomization argument involves treating the randomization indicators as random variables, and the observed outcomes as fixed. While the approach provides an attractive justification for using the repeated measures ANOVA in the randomized experiment, it may not be justifiable in the more general setting. In addition the justification relies on a linear model. Note that with constant variance and covariance, we can formulate a repeated measures analysis using mixed effects models with a random subject effect  $b_i$  (where  $\text{Var}(b_i) = \sigma_b^2$ ) and independent errors  $\epsilon_{ij}$  (where  $\text{Var}(\epsilon_{ij}) = \sigma_\epsilon^2$ ). Then  $\text{Var}(\hat{\delta}) = 2 \sigma_\epsilon^2/N$  and  $\text{Var}(\hat{\gamma}) = 2 (\sigma_\epsilon^2 + \sigma_b^2)/N$  because  $b_i$  drops out of any within-subject contrast.

Finally, there are many variations on repeated measures designs which are difficult to handle with the classical analytic approaches. For example, missing data lead to unbalanced designs with a variable number of observations on subjects. In addition the outcomes may be count or binary data, neither of which can be handled by classical ANOVA techniques. Using generalized linear models for correlated outcomes enables us to handle these variations in a unified way.

## 21.3 MULTIPLE SOURCE DATA

Multiple source data usually arise in the context of epidemiological studies, where outcomes and/or risk factors may be difficult to measure. For example, Fitzmaurice et al. (1995) discuss a study involving child psychopathology, which used standardized questionnaires administered to both the child's mother and teacher. Responses to these questionnaires were used to construct, for each informant, dichotomous indicators of whether psychopathology was present in the child. Other informants sometimes used in this setting include the child, the father, and peers. Because informants interact in different settings with the child, the information from different informants reflects different perspectives. However, the primary interest of the study was not to compare informants, but rather to study the effects of covariates on child psychopathology, broadly defined. By including informant as a variable in the correlated data model, and its interaction with other covariates of interest, one can also examine how informants differ overall, and whether covariate effects differ by informant.

In this section we use the term “multiple source” to encompass data that are simultaneously obtained from multiple informants or raters (e.g., self-reports, family members, health care providers, administrators) or via different/parallel instruments or methods (e.g., symptom rating scales, standardized diagnostic interviews, or clinical diagnoses). For example, in psychiatric studies of children, the child's parent is routinely used as a proxy data source; other sources (e.g., self-report, peers, teachers, clinicians, or trained observers) may also be employed, depending on the child's age and the nature of psychopathology under study. Multiple source data have become increasingly common in hospital-based and outpatient-based assessments of the effectiveness of treatments. For example, evaluations of managed care programs for the U.S. Medicaid population require analysis of multiple sources of information, including patient satisfaction with health care, treatment utilization, and appropriate care. Other areas where multiple reports arise include studies of severe mental illness, such as schizophrenia or Alzheimer's disease, where the affected subject is often unable to provide self-report data; family history studies, where many relatives are interviewed about the status of the proband and other family members; and behavioral studies of alcohol/drug use or of eating disorders, where information is obtained from the subject, as well as family members or other sources.

Historically there has been little consensus as to how to analyze multiple source data. Sometimes investigators conduct completely separate univariate analyses for each source. Alternatively the multiple source measures may be combined to make a single outcome. For measured responses, investigators may take a mean and for dichotomous responses the “and” or the “or” rules are often used. In the “and” rule binary source data are considered to be positive if *all* of the source data are positive, and negative otherwise; in the “or” rule binary source data are considered to be positive if *any* of the source data are positive, and negative otherwise. All of these ad hoc strategies usually require discarding data when any sources are missing; the separate analysis strategy makes it difficult to compare the results for the two (or more) sources, while the analysis using a combined response can obscure interesting differences. Using correlated data models similar to those discussed in previous chapters allows one to directly compare source effects and to handle missing data in a unified framework.

## 21.4 CASE STUDY 1: REPEATED MEASURES EXPERIMENT

In this section we analyze data from a randomized crossover design to illustrate the use of mixed effects models in the repeated measures setting. The study was designed to compare two active drugs and placebo for relief of tension headache. The two analgesic drugs were identical in their active ingredients except that one contained caffeine. The primary comparison of interest was between the two active drugs; the placebo was included for regulatory purposes. What makes the analysis non-standard is that there were three treatments, but only two periods; that is, each subject received only two of the three treatments in a random order. With three treatments and two periods, there are six possible treatment sequences, AB, BA, AP, PA, BP, PB, where A, B, and P denote the two active drugs and placebo. If each sequence is assigned an equal number of subjects, then we have what is known as a balanced incomplete block design. It is balanced because all possible sequences are equally represented, but incomplete because each subject is “missing” a third treatment. In our example the AB and BA sequences were assigned three times as many subjects as the remaining four because of the interest in the A versus B comparison. The descriptive statistics for one measure of pain relief used in the crossover study of analgesics are given in [Table 21.1](#). There were actually two headaches treated within each period, but that feature of the data is ignored here and we use the average of the two measures of pain relief. A few observations were missing, but we have access only to subjects with complete data.

**Table 21.1** Descriptive statistics (means, standard deviations, and correlation), by sequence, for total pain relief for headache in periods 1 and 2.

Sequence	N	Period 1		Period 2		$\rho$
		Mean	SD	Mean	SD	
AB	126	10.196	3.347	9.153	3.429	0.20
BA	127	9.581	3.881	10.791	3.530	0.30
AP	43	10.477	3.546	7.273	4.451	0.31
PA	43	8.366	3.777	10.855	3.204	0.47
BP	42	10.333	3.306	8.357	3.944	0.72
PB	42	7.464	4.265	9.911	4.183	0.67

Note:  $\rho$  is the correlation between pain relief scores in periods 1 and 2. The estimate of the common correlation (ignoring sequence) is  $\hat{\rho} = 0.38$ .

An important issue in crossover designs is the possibility of carryover effects. The presence of a carryover effect means that the treatment taken in the first period may influence the treatment effect in the second period, that is, the drug in the first period carries over. In our analysis we will show the results of fitting two different models, one with only treatment and period effects included in the model and one with the carryover effects included. We display the model that includes carryover effects by giving the expected value of the response for each sequence and each period. Letting  $Y_{ij}$  denote the treatment response for the  $i^{th}$  subject in the  $j^{th}$  period, we write

$$(21.1) \quad E(Y_{ij}) = \beta_1 + \beta_2 X_{2ij} + \beta_3 X_{3ij} + \beta_4 X_{4ij} + \beta_5 X_{5ij} + \beta_6 X_{6ij},$$

where  $X_{1ij} = 1$  for all  $i$  and  $j$ ,  $X_{2ij}$  and  $X_{3ij}$  are dummy variable indicators of the main effects of treatment A and B (P is the reference condition),  $X_{4ij}$  is a dummy variable indicator of period 2 (so for all  $i$ ,  $X_{4i1} = 0$  and  $X_{4i2} = 1$ ) and  $X_{5ij}$  and  $X_{6ij}$  are dummy variable indicators of the carryover effects for A and B, respectively. For subjects assigned to treatment A in the  $j^{th}$  period,  $X_{2ij} = 1$ , and 0 otherwise; for those assigned B in the  $j^{th}$  period,  $X_{3ij} = 1$ , and 0 otherwise ( $X_{2ij} = X_{3ij} = 0$  for subjects assigned to P in the  $j^{th}$  period). For subjects in sequences where A is taken first  $X_{5i2} = 1$ , and  $X_{6i2} = 1$  if the first treatment taken is B; for all other values of  $i$ ,  $X_{5i2}$  and  $X_{6i2}$  are zero. Likewise, for

all  $i$ ,  $X_{5i1} = X_{6i1} = 0$  since, by definition, there are no carryover effects in the first period. In [Table 21.2](#) we summarize the interpretations of  $\beta_1, \dots, \beta_6$  in terms of the mean response for each sequence during each period.

**Table 21.2** Regression parameters for the mean total pain relief for headache, by sequence and period.

Sequence	Period 1	Period 2
AB	$\beta_1 + \beta_2$	$\beta_1 + \beta_3 + \beta_4 + \beta_5$
BA	$\beta_1 + \beta_3$	$\beta_1 + \beta_2 + \beta_4 + \beta_6$
AP	$\beta_1 + \beta_2$	$\beta_1 + \beta_4 + \beta_5$
PA	$\beta_1$	$\beta_1 + \beta_2 + \beta_4$
BP	$\beta_1 + \beta_3$	$\beta_1 + \beta_4 + \beta_6$
PB	$\beta_1$	$\beta_1 + \beta_3 + \beta_4$

This model includes effects of treatment ( $\beta_2$  and  $\beta_3$ ), period ( $\beta_4$ ), and carryover ( $\beta_5$  and  $\beta_6$ ). The active drug comparison is given by  $\beta_2 - \beta_3$ . In our model the carryover effects are assumed to depend only on the treatments taken in period one. That is, the model assumes that the carryover of one active treatment to the other active treatment is the same as the carryover of that active treatment to the placebo. More complicated models can be used (e.g., see Laird, Skinner, and Kenward, 1992) but are not considered here.

Since subjects are randomized, we will assume  $\text{Var}(Y_{i1}) = \text{Var}(Y_{i2})$ . A more general covariance structure could easily be accommodated by allowing the variance to depend on period, but we do not do so here. Using the compound symmetry assumption allows us to analyze the data using mixed effects models with a random subject effect. Handling missing responses simply involves removing rows from the design matrix and the response vector (as discussed in Section 5.3).

The balanced incomplete block design has two types of information about treatment. One source of information comes from within-subject contrasts; that is, in the absence of carryover effects,  $Y_{i2} - Y_{i1}$  is a treatment contrast with variance  $2\sigma^2(1 - \rho)$ . There is also information about treatment contrasts from between-subject comparisons; that is, in the absence of carryover effects,  $Y_{ij} - Y_{i'j}$  estimates a treatment contrast for two subjects with different treatments in period  $j$ . Now, however, the variance of the contrast is  $2\sigma^2$ . When there are carryover effects in this design, information about treatment effects can be found in both within-subject and between-subject contrasts. If  $\rho$  is modest ( $\rho < 0.4$ ), it is often suggested that the between-subject information be ignored, and the analysis use only the within-subject contrasts. The rationale behind this approach is that the within-subject contrast yields a simple structure with no need to estimate variance and covariance components for  $Y_{i1}$  and  $Y_{i2}$ , while using all the information requires estimating the between- and within-subject error variance. Of note,  $\hat{\rho} \approx 0.4$  for the pain relief data in our example (see [Table 21.1](#)).

The within-subject analysis is especially easy to do with just two periods by subtracting  $Y_{i1}$  from  $Y_{i2}$ , and analyzing the difference  $d_i = Y_{i2} - Y_{i1}$ :

$$(21.2) E(d_i) = \beta_2(X_{2i2} - X_{2i1}) + \beta_3(X_{3i2} - X_{3i1}) + \beta_4 + \beta_5X_{5i2} + \beta_6X_{6i2}.$$

Note that  $\beta_1$  has vanished from the model, and  $\beta_4$  acts as the constant (or intercept) in the model, since  $X_{4i1} = 0$  and  $X_{4i2} = 1$  for all  $i$ . The parameters  $\beta_2$  and  $\beta_3$  remain the main effects of treatment and  $\beta_5$  and  $\beta_6$  remain the carryover effects. In our analysis the main focus is on the active drug comparison,  $\beta_2 - \beta_3$ . It is a special feature of this design that carryover effects can be estimated from within-subject contrasts (Koch et al., 1989), but as we will show, most of the information about the carryover effects comes from the between-subject information. This issue is similar to that raised in Chapter 9 (see Section 9.5) in the discussion of cross-sectional and longitudinal information. In the context of a balanced incomplete randomized trial, or a complete randomized trial with carryover

effects, the issue of bias does not arise, only that of efficiency.

We will illustrate the analysis using both the simple regression on the differences ( $d_i$ ) and a mixed effects model analysis on the full data ( $Y_{i1}$  and  $Y_{i2}$ ), with subject as a random effect ( $b_i$ ). [Table 21.3](#) illustrates the conventional wisdom that the within-subject (differences) analysis is highly efficient for treatment effects when carryover is absent. (Compare the two rows of [Table 21.3](#) for the treatment effect estimates assuming no carryover effects.) However, if carryover effects are present the within-subject analysis is very inefficient for both the main effect of treatment and carryover effects. (Compare the two rows of [Table 21.3](#) for the treatment and carryover effects estimates.) In this case investigators basing their conclusion on the within-subject (differences) analysis would be led to conclude erroneously that the active drug comparison is not confounded with carryover ( $Z = -1.19/1.26 = -0.94, p > 0.30$ ) and to use the results of the model without carryover to obtain a significant difference between the active treatments ( $Z = 1.06/0.24 = 4.42, p < 0.0001$ ). That is, the inefficient analysis of carryover effects based on differences leads to the erroneous conclusion that the treatment comparison does not require any adjustment for carryover effects. Using the mixed effects model analysis, however, we come to quite a different conclusion; there is a substantial carryover effect ( $Z = -1.06/0.52 = -2.05, p < 0.05$ ), but no evidence of a statistically significant treatment effect ( $Z = 0.55/0.32 = 1.71, p > 0.05$ ) when the carryover effects are included in the model and adjusted for in the analysis. The widespread availability of software to implement a mixed effects model analysis makes it relatively easy to capture the “between-subject information” even in complex repeated measures designs.

**Table 21.3** Results of comparison of the two active treatments (estimate  $\pm$  SE) based on analysis of differences and mixed effects model analysis, with and without carryover effects.

Method of Analysis	No Carryover		With Carryover	
	Treatment	$(\beta_2 - \beta_3)$	Treatment	Carryover
Differences	$1.06 \pm 0.24$		$0.46 \pm 0.67$	$-1.19 \pm 1.26$
Mixed Effects Model <sup>a</sup>	$1.02 \pm 0.23$		$0.55 \pm 0.32$	$-1.06 \pm 0.52$

<sup>a</sup>Linear mixed effects model with random subject effect.

## 21.5 CASE STUDY 2: MULTIPLE SOURCE DATA

Data for this example come from two surveys of children's mental health (Zahner et al., 1992, 1993). A standardized measure of childhood psychopathology was used both by parents (Child Behavior Checklist, CBC) and teachers (Teacher Report Forms, TRF) to assess children in the study. We use here the externalizing scale, which assesses delinquent and aggressive behavior. The scale has been dichotomized at the cut point for borderline/clinical psychopathology. The cut points are normed separately for males and females; thus we expect to see small gender effects in these data. Because of the multiple levels of permissions and reporting, a substantial number of children were missing the TRF. Our analysis is based on 1428 children who had both parent and teacher responses, and an additional 1073 children with only a parent response; a total of 2501 children participated in the study. In this example the two sources or respondents are the children's parents and teachers; in the psychiatric literature these sources are often referred to as "informants."

The objective of the analysis is to study the influence of several explanatory variables on the prevalence of externalizing behavior in these children. For simplicity we limit our analysis to single-parent status (coded 1: single, 0: otherwise) and child's physical health problems (coded 1: fair to bad health, 0: good health). In addition we are interested in describing the level of association between the two respondents, and determining whether the effects of the covariates depend on informant. We will use standard regression models to address these issues, but because we have two different measures of the outcome, we will use correlated data models, one for each outcome. Since externalizing behavior is dichotomous, we use logistic models for the regressions.

The basic approach uses two separate regression models, one with the CBC as an outcome and one with the TRF as an outcome. Both models have the same set of covariates (here single-parent status and physical health problems), but the coefficients may differ for the different sources. In addition we have an "informant" indicator variable that identifies the source of the response. To motivate the approach, we first fit two completely separate logistic regressions each with the full complement of covariates, one for each informant outcome. Let  $\mu^P$  and  $\mu^T$  denote the probability of a positive response on externalizing behavior as measured by parents and teachers, respectively. Then the two regression models are

$$\text{logit}(\mu_i^P) = X_i' \beta^P$$

and

$$\text{logit}(\mu_i^T) = X_i' \beta^T,$$

where  $X_i$  is a  $p \times 1$  vector of covariates for the  $i^{th}$  subject and

$$E(Y_i^P | X_i) = \mu_i^P \text{ and } E(Y_i^T | X_i) = \mu_i^T.$$

The first logistic regression model was fit with the parent response as outcome using all 2501 children, and the second was fit with the teacher response using only the 1428 children with a teacher response. The results are displayed in [Table 21.4](#).

**Table 21.4** Results of fitting separate logistic regressions to data on externalizing behavior from each source.

Informant	N	Intercept <sup>a</sup>	Single Parent <sup>a</sup>	Child Health <sup>a</sup>
Parent	2501	$-2.156 \pm 0.092$	$0.620 \pm 0.124$	$0.600 \pm 0.113$
Teacher	1428	$-1.694 \pm 0.105$	$0.655 \pm 0.157$	$0.175 \pm 0.135$

<sup>a</sup>Estimated coefficient  $\pm$  standard error.

The estimated coefficients for single parent status are similar and statistically significant (at the 0.05 level) for each informant report; the estimated coefficients for child health problems are rather different, and significant only for the parent report. In both cases standard errors for the parent

response are smaller, reflecting the larger sample size. Fitting these two logistic regression models separately does not allow us to formally quantify the differences in  $\beta^P$  and  $\beta^T$  because the estimated regression parameters are correlated.

We now show how to fit these two regression models simultaneously using multivariate methods that take the association between the responses into account. To begin, we rewrite the two separate models as a set of bivariate models with a common regression coefficient  $\beta$ , which will have dimension six (or  $2p$  in general). First, we simply change notation as follows. Let  $Y_i^P = Y_{i1}$ ,  $Y_i^T = Y_{i2}$ ,  $\mu_i^P = \mu_{i1}$ ,  $\mu_i^T = \mu_{i2}$ ,  $\beta' = (\beta^P, \beta^T)$ , and let  $Z_i$  be an  $n_i \times 6$  matrix where  $n_i$  is the number of informants available for the  $i^{th}$  child. For a child with both informants ( $n_i = 2$ ), the first row of  $Z_i$  is given by

$$Z'_{i1} = (X'_i, 0, 0, 0)$$

and the second row of  $Z_i$  simply interchanges  $X'_i$  with  $(0, 0, 0)$ ,

$$Z'_{i2} = (0, 0, 0, X'_i).$$

If an informant is missing ( $n_i = 1$ ), we delete the row corresponding to that informant (i.e., delete  $Z_{i2}$  if the teacher does not give a report for the  $i^{th}$  child). We now write

$$E(Y_{ij}) = \mu_{ij}, \quad i = 1, \dots, N \text{ and } j = 1, \dots, n_i,$$

and specify the following bivariate model,

$$(21.3) \text{ logit}(\mu_{ij}) = Z'_{ij}\beta, \quad j = 1, 2.$$

This is exactly the mean model for a marginal model (see Chapters 12 and 13), with a special structure for the design matrix, here labeled  $Z_i$  instead of the usual  $X_i$ . To complete the marginal model, all we need is a specification of  $\text{Cov}(Y_{i1}, Y_{i2})$ , using one of the approaches discussed in Chapters 12 and 13. For the analysis here we use the odds ratio to measure the association. Given the model for the mean and  $\text{Cov}(Y_{j1}, Y_{j2})$ , GEE methods can be used to estimate  $\beta$ . If we specify the entire joint distribution for the  $Y_{ij}$ 's, we can use maximum likelihood estimation (Fitzmaurice et al., 1995). We choose to present a GEE analysis since it is easier to implement. Note that, because of the way we defined  $Z_i$  and  $\beta$ , the first three components of  $\beta$  correspond to  $\beta^P$  and the second three correspond to  $\beta^T$ . The variance-covariance matrix of  $\hat{\beta}$  is now  $6 \times 6$ ; the  $3 \times 3$  diagonal blocks are the variance-covariance matrices of  $\hat{\beta}^P$  and  $\hat{\beta}^T$ , while the off-diagonal  $3 \times 3$  block gives the covariance between the two.

[Table 21.5](#) shows the results of fitting the regression model, using the log odds ratio to model association between parent and teacher response. The estimates and standard errors for  $\beta^P$  are nearly unchanged, as we might expect, but those for  $\beta^T$  are different, reflecting the fact that many children were missing the teacher response. The parent response provides some information about teacher response because of the relatively high association between parent and teacher responses (estimated odds ratio is 4.75, with a 95% confidence interval of 3.52 to 6.39). For both  $\hat{\beta}^T$  and  $\hat{\beta}^P$  the standard errors are slightly smaller than with separate logistic regressions, but only very slightly so for the  $\hat{\beta}^P$ .

[Table 21.5](#) Results of fitting two regression models simultaneously to externalizing behavior data on 2501 children using GEE method.

Informant	Intercept <sup>a</sup>	Single Parent <sup>a</sup>	Child Health <sup>a</sup>
Parent	$-2.154 \pm 0.091$	$0.616 \pm 0.124$	$0.598 \pm 0.113$
Teacher	$-1.683 \pm 0.104$	$0.602 \pm 0.155$	$0.146 \pm 0.135$

<sup>a</sup>Estimated coefficient  $\pm$  empirical standard error.

If we use a “working independence” model for  $\text{Cov}(Y_{i1}, Y_{i2})$ , setting the log odds ratio for the association between parent and teacher response to zero, then the results (not shown) of a GEE analysis yield exactly the same result as separate regressions for  $\hat{\beta}$  (i.e., GEE estimates of  $\beta$  are identical to those reported in [Table 21.4](#)). The model-based standard errors are also identical to

those reported in [Table 21.4](#) and the  $3 \times 3$  off-diagonal block of the covariance matrix for  $\hat{\beta}$  is zero. The empirical (or “sandwich”) standard errors are very similar to the model-based standard errors. However, they differ slightly because the  $3 \times 3$  off-diagonal block of the covariance matrix for  $\hat{\beta}$  is estimated as zero by the model-based variance estimator, whereas the empirical variance estimator correctly estimates the covariance between  $\hat{\beta}^P$  and  $\hat{\beta}^T$ . That is, the empirical standard errors account for the correlations (ranging from approximately 0.15 to 0.25) between components of  $\hat{\beta}^P$  and  $\hat{\beta}^T$ .

The model given by [\(21.3\)](#) is a very general model; its advantages over the separate regressions are that:

1. we can test whether  $\beta_k^P = \beta_k^T$  for the  $k^{th}$  covariate (or for the whole vector, test  $\beta^P = \beta^T$ , using  $\text{Cov}(\hat{\beta})$  provided by the GEE analysis);
2. we can use all available data; and
3. it provides a measure of association between the two informants.

With a large number of covariates, we will usually want to fit simpler models. The way we have defined  $\beta$  and  $Z_i$  in model [\(21.3\)](#) means that the first  $p$  components of  $\beta$  correspond to  $\beta^P$  and the second  $p$  correspond to  $\beta^T$ . To formulate simpler models, we need to create a dichotomous indicator variable of informant.

To illustrate, we introduce a dichotomous variable ( $X_2$ ) which is 1 if the informant is the parent, and 0 if the informant is the teacher. Denoting single-parent status and child health problems by  $X_3$  and  $X_4$ , consider a model of the form

$$(21.4) \text{logit}(\mu_{ij}) = \beta_1 + \beta_2 X_{2ij} + \beta_3 X_{3ij} + \beta_4 X_{4ij},$$

where  $X_{1ij} = 1$  for all  $i$  and  $j$ . This model specifies that the effects of single-parent status (measured by  $\beta_3$ ) and child’s physical health problems (measured by  $\beta_4$ ) are the same regardless of the source of the information, but the mean level may be higher or lower (measured by  $\beta_2$ ) depending on informant. A positive  $\beta_2$  suggests that a positive rating ( $Y_{ij} = 1$ ), here denoting externalizing behaviors in the borderline/clinical range, is more likely from a parent’s report than from a teacher’s, holding single parent status and physical health problems constant. Notice that only informant is a within-subject variable; that is,  $X_{2i1} \neq X_{2i2}$  while  $X_3$  and  $X_4$  are both between-subject variables. Forcing the coefficients of single-parent status to be equal seems reasonable in view of the results presented in [Tables 21.4](#) and [21.5](#), but not for child health problems.

We construct a model which allows the effect of physical health problems to depend on informant by simply adding the interaction:

$$(21.5) \text{logit}(\mu_{ij}) = \beta_1 + \beta_2 X_{2ij} + \beta_3 X_{3ij} + \beta_4 X_{4ij} + \beta_5 X_{5ij},$$

where  $X_{5ij} = X_{2ij} X_{4ij}$ . Fitting a model that includes the interaction allows the probability of a positive rating to depend on informant, single-parent status, and child health problems, and allows the effect of child health problems to differ by informant. The results of this model are displayed in [Table 21.6](#). The estimated odds ratios change very little from the model given by [\(21.3\)](#). The interpretation of the results in [Table 21.6](#) becomes more transparent if the model given by [\(21.5\)](#) is written separately for the parent and teacher informant:

$$\begin{aligned} \text{logit}(\mu_{i1}) &= \text{logit}(\mu_i^P) = (\beta_1 + \beta_2) + \beta_3 X_{3ij} + (\beta_4 + \beta_5) X_{4ij}; \\ \text{logit}(\mu_{i2}) &= \text{logit}(\mu_i^T) = \beta_1 + \beta_3 X_{3ij} + \beta_4 X_{4ij}. \end{aligned}$$

**Table 21.6** Results of fitting regression models, with common or shared parameters, simultaneously to data on externalizing behavior using GEE method.

Variable	Estimate	SE <sup>a</sup>	Z
Intercept	-1.685	0.100	-16.85
Informant	-0.467	0.118	-3.96
Single Parent	0.611	0.108	5.68
Child Health	0.146	0.135	1.08
Informant × Child Health	0.452	0.157	2.87

<sup>a</sup>Empirical standard error.

Because we code informant as 1 for parent response,  $\hat{\beta}_1$  and  $\hat{\beta}_4$  can be regarded as estimated teacher parameters for the intercept and child physical health problems. For single-parent status,  $\hat{\beta}_3$  is the common coefficient for both informants; notice that its standard error is considerably smaller than the two corresponding standard errors for model (21.3) reported in [Table 21.5](#). Finally, from  $\beta_5$  we have the difference in effects of child health problems as estimated by parent and teacher evaluations. Comparison of  $\hat{\beta}_5$  to its standard error provides a formal test of the null hypothesis that informant does not matter in evaluating the effect of child health problems. The rejection of this hypothesis may reflect the fact that the rating of child health problems was given by the parent, and the teacher may lack information about child's physical health.

As a rule, if informant interactions are included for all the covariates, then the model is basically equivalent to fitting separate regressions. However, the estimated coefficients obtained from using GEE with a non-zero correlation will differ from those obtained by fitting separate regressions because of the non-linearity of the logistic regression model, even with complete data on each informant. If there are missing data, the coefficients may differ substantially. In our example the estimated effects for the parent respondents are very similar to those obtained from a separate regression on  $Y_{i1}$ , with all 2501 observations. The differences for the teacher respondents were more pronounced because of the substantial missing data, and because the two informant responses are highly associated. When we use GEE with a non-zero correlation, the coefficients for the teacher responses use the information in the parent response to provide some information for the missing teacher respondents.

When simpler models are fit (i.e., not all interactions with informant are present) we can expect to gain efficiency in the analysis for the common coefficients. This point is illustrated by comparing the standard errors of the coefficient for single-parent status in [Table 21.4](#) or 21.5 with [Table 21.6](#).

## 21.6 SUMMARY

This chapter has illustrated how two other types of studies, repeated measures and multiple sources, can be analyzed using methods for correlated data that are comparable to those used for longitudinal data analysis. This is certainly not an exhaustive list, as many study designs produce repeated measures and multiple source data, and their proper analysis requires linear or generalized linear models for correlated data.

Our repeated measures example is also an example of a crossover design, but there are many examples of simpler repeated measures designs with no period or carryover effects; in these cases there are often between-subject variables and the interaction of those with the within-subject variables will be of interest. For our example the response variable is continuous and a linear mixed model is appropriate, but in other settings one may wish to use a more general covariance structure or, when the response variable is discrete, a generalized linear model for correlated data.

Our multiple informant example used dichotomous reports and logistic regression models, but often multiple source data will be continuous responses. The general approach to constructing multivariate correlated regression models remains the same, but now maximum likelihood and GEE approaches are viable alternatives for the analysis. Likewise there may be more than two sources ( $n > 2$ ). The general approach to analyzing multiple source data can handle any number of informants or sources using  $n - 1$  (one less than the number of informants) dichotomous indicators of informant and their interactions with the covariates of interest.

## 21.7 FURTHER READING

Repeated measures designs are frequently encountered in applications and there is a very large literature on their design and analysis. Accessible descriptions of methods for analyzing repeated measures data can be found in the review articles by Keselman and Keselman (1984), Everitt (1995), and Omar *et al.* (1999).

Methods for analyzing crossover trials are discussed in the review articles by Matthews (1994) and Jones and Donev (1996). Finally, Fitzmaurice et al. (1995) and Goldwasser and Fitzmaurice (2001) discuss the use of regression models for analyzing multiple source data and present a substantive analysis of the multiple informant data from the Connecticut Child Surveys.

## Bibliographic Notes

A comprehensive description of the multitude of techniques available for analyzing repeated measures data can be found in the books by Crowder and Hand (1990), Lindsey (1999), and Davis (2002), and the references therein.

A discussion of split-plot designs can be found in Chapter 7 of Cochran and Cox (1957). A comprehensive description of the design and analysis of crossover trials can be found in the books by Jones and Kenward (1989) and Senn (2002), and the references therein.

# *Chapter 22*

## *Multilevel Models*

### **22.1 INTRODUCTION**

In Parts I through IV the major focus has been on the analysis of longitudinal data. A distinctive feature of longitudinal data is that they are *clustered*. In longitudinal studies the cluster is composed of the repeated measurements obtained from a single individual at different occasions. There are, however, many studies in the health sciences that are not longitudinal but that give rise to data that are clustered or cluster-correlated. For example, clustered data commonly arise when intact groups are randomized to health interventions in so-called cluster-randomized trials or when naturally occurring groups in the population are randomly sampled. In addition there can be more than a single level of clustering in the data. The term *multilevel data* (or *hierarchical data*) encompasses all of these cases. The distinctive feature of multilevel data is that measurements on units within a cluster are more similar than measurements on units in different clusters. The clustering can be expressed in terms of correlation among the measurements on units within the same cluster and this correlation must be appropriately accounted for in the analysis.

Because longitudinal data are a special case of multilevel data, with only a single level of clustering and a natural ordering of the measurements within a cluster, this chapter provides a description of regression models for multilevel data, more broadly defined. One of our goals is to demonstrate that many of the methods for the analysis of longitudinal data considered in earlier chapters are, more or less, special cases of more general regression methods for multilevel data. The overview of multilevel models presented in this chapter provides a basic introduction to a general methodology for analyzing the wide range of clustered data that commonly arise in studies in the biomedical and health sciences.

## 22.2 MULTILEVEL DATA

Multilevel data arise when there is a hierarchical or clustered structure to the data. Data of this kind frequently arise in the health sciences, since individuals can be grouped in so many different ways. For example, in studies of health services and outcomes, assessments of quality of care are often obtained from patients who are nested within different clinics. Such data can be regarded as multilevel, with patients referred to as the level 1 units and clinics the level 2 units. In this example there are two levels in the data hierarchy and, by convention, the lowest level of the hierarchy is referred to as level 1. The term “level,” as used in this context, signifies the position of a unit of observation within a hierarchy.

Broadly speaking, the clustering in multilevel data can be a consequence of the study design or due to a naturally occurring hierarchy in the target population, or sometimes due to both. An example of a naturally occurring two-level data hierarchy arises in developmental toxicity studies. In a typical developmental toxicity experiment, pregnant mice or rats (dams) are assigned to increasing doses of a chemical or a test substance over the period of major organogenesis (when organ systems are developing in a growing fetus). Following sacrifice, each fetus in the litter is weighed (a continuous response) and examined for evidence of malformations (a binary response, present or absent). Data collected in developmental toxicity experiments are clustered (i.e., the litter is the cluster), with observations on the fetuses (level 1 units) nested within dams/litters (level 2 units). Two-level data also arise in family studies designed to assess the association or “aggregation” of disease (or markers of disease development) among relatives. In family studies the goal is to determine whether the presence of disease in a family member is associated with increased risk of disease for relatives. The associations among siblings and between parent-child pairs are of main interest because their relative magnitudes can be used to provide indirect evidence of genetic risk due to the sharing of the same genes. Data from studies of nuclear families are clustered, with observations on the mother, father, and children (level 1 units) nested within families (level 2 units).

Other common examples of naturally occurring clusters in the population are households, hospital wards, clinics, medical practices, neighborhoods, and schools. Furthermore naturally occurring hierarchical data structures can have more than two levels. For example, observations may be obtained on patients nested within clinics, which in turn are nested within different geographical regions of the country. Another example of a naturally occurring data hierarchy is when observations are obtained on children nested within classrooms, nested within schools. In both of these examples there are three levels in the data hierarchy. In principle, there can be many levels in the data hierarchy.

Alternatively, the hierarchical data structure can be a consequence of the study design. For example, the U.S. National Health and Nutrition Examination Survey (NHANES) uses a multi-stage sampling design to produce information on nutrition and health status. The target population is the total U.S. civilian non-institutionalized population, 2 months of age or over. Because it is not practical to obtain a simple random sample of the U.S. population, complex sampling methods are commonly used. For example, NHANES III, conducted in 1988 to 1994, used the following multi-stage sampling design (National Center for Health Statistics, 1992, 1994). In the first stage, so-called primary sampling units (PSUs) were defined based on counties or combined counties in the United States. A first-stage random sample of PSUs was selected from these geographical regions. In the second-stage sampling, within each of the selected PSUs, a random sample of area segments consisting of census blocks was selected. In the third stage, within each of the selected area segments, a random sample of households was selected. Finally, in the fourth stage, eligible persons were randomly selected within households. The resulting data can be regarded as hierarchical, with individuals being the level 1 units, households the level 2 units, area segments the level 3 units, and counties the level 4 units.

Additional examples of study designs that produce multilevel data structures include cluster-randomized clinical trials, repeated measures experiments, and longitudinal studies. In a cluster-randomized trial, groups of individuals, rather than the individual subjects, are randomized to

different treatments or health interventions. For example, the Promotion of Breastfeeding Intervention Trial (PROBIT) was designed to determine whether efforts to promote breastfeeding have any impact on the duration and exclusivity of breastfeeding (Kramer et al., 2001). In this trial maternity clinics, rather than the mothers, were randomized to either the intervention or control (standard care). The mothers were followed-up for one year after the birth of their infants and the effectiveness of the health intervention was assessed by the responses of mothers in each treatment group. When regarded as multilevel data, the level 1 units are the mothers and the level 2 units are the maternity clinics. Of note, the main covariate of interest, denoting the assignment to intervention or control, is defined at level 2. Longitudinal studies are another common example where the study design produces data with a two-level structure. In a longitudinal study the clusters are composed of the repeated measurements obtained from a single individual at different occasions. When longitudinal data are regarded as multilevel data, the level 1 units are the repeated occasions of measurement and the level 2 units are the subjects.

Finally, the clustering in multilevel data can be due to both the design of the study and naturally occurring hierarchies in the target population. For example, clinical trials are often conducted in many centers to ensure sufficient numbers of patients and/or to assess the effectiveness of the treatment in different settings. These studies are referred to as multi-center trials. Observations from a multi-center longitudinal clinical trial can be regarded as multilevel data having 3 levels, with repeated measurement occasions (level 1 units) nested within subjects (level 2 units) nested within centers (level 3 units).

Although we have distinguished between clustering that occurs naturally and clustering due to study design, the consequence of clustering at different levels is the same: units that are grouped at any level are likely to respond more similarly. For example, two patients selected at random from the same clinic are expected to respond more similarly than two patients randomly selected from different clinics. In general, the degree of clustering can be expressed in terms of correlation among the observations on units within the same level. Statistical models for multilevel data must account for the intra-cluster correlation at each level; failure to do so can result in misleading inferences.

## 22.3 MULTILEVEL LINEAR MODELS

In this section we discuss linear models for multilevel data. The dominant approaches to multilevel modeling have the same basis: clustering in the data is accounted for via the introduction of random effects at different levels in the hierarchy. Multilevel linear models can be regarded as extensions of the linear mixed effects models described in Chapter 8, which allow random effects to be incorporated at more than one level. In addition to accounting for clustering in the data, multilevel models permit estimation of the effects of covariates, measured at any of the levels of the hierarchy, on the outcome.

In a multilevel model the response is obtained on the lowest level (or level 1) units, but covariate information can be measured at any level. Combining covariates measured at different levels of the hierarchy within a single regression model is central to multilevel modeling. For example, multilevel models can determine and disentangle the relative importance of patient-level, clinic-level, and region-level factors on quality of care. In general, multilevel models can be used to make inferences about the population of units at any level of the hierarchy and to discern how variation in the outcome at different levels depends on covariates. In this section we present an overview of multilevel linear models for a continuous outcome. We begin with a discussion of models for two-level data. The models generalize in a natural way when there is additional clustering in the data (e.g., three-level and higher-level data). The major focus of this section is on the specification of multilevel models; estimation is mentioned but not emphasized.

## 22.3.1 Two-Level Linear Models

Before describing models for two-level data we need to introduce some notation. For a two-level data structure, let  $i$  index level 1 units and  $j$  index level 2 units. We assume that there are  $n_2$  units at level 2 in the sample. Each of these clusters (for  $j = 1, \dots, n_2$ ) is composed of  $n_{1j}$  level 1 units. For example, consider a multi-center clinical trial comparing two treatments (active drug versus placebo) conducted in 20 medical clinics. Patients are enrolled from each clinic and randomly assigned to one of the two treatment conditions. In this example, clinics are the level 2 units ( $j = 1, \dots, 20$ ) and patients are the level 1 units ( $i = 1, \dots, n_{1j}$ ), where  $n_{1j}$  is the number of patients enrolled in the study from the  $j^{th}$  clinic (and  $n_2 = 20$  is the number of clinics). Alternatively, consider two-level data arising from a longitudinal study where 150 subjects are measured at four occasions. In this example, subjects are the level 2 units ( $j = 1, \dots, 150$ ), and measurement occasions are the level 1 units ( $i = 1, \dots, 4$ ), with  $n_2 = 150$  level 2 units, and  $n_{1j} = 4$  level 1 units (within each level 2 unit). For the latter example the alert reader will have noticed that the indices  $i$  and  $j$  have now been reversed from their use in earlier chapters; here we have adopted the usual convention in much of the multilevel modeling literature of letting  $i$  denote level 1 units,  $j$  denote level 2 units, and so on. We must caution the reader that some of the literature on multilevel modeling reverses this notation (and/or occasionally reverses the ordering of the levels).

Let  $Y_{ij}$  denote the response on the  $i^{th}$  level 1 unit within the  $j^{th}$  level 2 cluster. For example,  $Y_{ij}$  might denote the primary outcome for the  $i^{th}$  patient in the  $j^{th}$  clinic. Associated with each  $Y_{ij}$  is a  $1 \times p$  (row) vector of covariates,  $X_{ij}$ . These can include covariates defined at each of the two levels and can also include “compositional” covariates, so called because they are formed by aggregating values over lower level units. For example, severity of disease defines a patient-level (or level 1) covariate. However, a “compositional” covariate at the clinic level can be formed by taking the average disease severity for all patients within each clinic.

Consider the following linear model relating the mean response to the covariates:

$$(22.1) \quad E(Y_{ij}|X_{ij}) = X_{ij}\beta.$$

For example, in a multi-center clinical trial, a simple model for the mean response is given by

$$E(Y_{ij}) = \beta_1 + \beta_2 \text{Group}_{ij},$$

where  $\text{Group}_{ij}$  denotes the treatment assignment for the  $i^{th}$  patient in the  $j^{th}$  clinic, with  $\text{Group}_{ij} = 1$  for active drug and  $\text{Group}_{ij} = 0$  for placebo. The model given by (22.1) specifies how the mean response depends on covariates, where the covariates can be defined at level 2 and/or level 1. A multilevel model accounts for the variability in  $Y_{ij}$ , around its mean, by allowing for random variation across both level 1 and level 2 units. In particular, a multilevel model for  $Y_{ij}$  assumes there is random variation across level 1 units and random variation in a subset of the regression parameters across level 2 units. The two-level linear model for  $Y_{ij}$  is given by

$$(22.2) \quad Y_{ij} = X_{ij}\beta + Z_{ij}b_j + \epsilon_{ij},$$

where  $Z_{ij}$  is a design matrix for the random effects at level 2, formed from a subset of the appropriate columns of  $X_{ij}$ . The random effects,  $b_j$ , vary across level 2 units but, for a given level 2 unit, are constant for all level 1 units. These random effects are assumed to be independent across level 2 units, with mean zero and covariance,  $\text{Cov}(b_j) = G$ . The level 1 random components,  $\epsilon_{ij}$ , are also assumed to be independent across level 1 units, with mean zero and variance,  $\text{Var}(\epsilon_{ij}) = \sigma^2$ . In addition the  $\epsilon_{ij}$ 's are assumed to be independent of the  $b_j$ 's, with  $\text{Cov}(\epsilon_{ij}, b_j) = 0$ . That is, the level 1 units are assumed to be conditionally independent given the level 2 random effects (and the covariates).

The regression parameters,  $\beta$ , are the fixed effects and describe the effects of covariates on the mean response

$$E(Y_{ij}) = X_{ij}\beta,$$

where the mean response is averaged over both level 1 and level 2 units. The two-level model given by (22.2) also describes the effects of covariates on the conditional mean response

$$E(Y_{ij}|b_j) = X_{ij}\beta + Z_{ij}b_j,$$

where the response is averaged over level 1 units only.

Let us return to the multi-center clinical trial example introduced earlier. A simple two-level model for the data is given by

$$Y_{ij} = \beta_1 + \beta_2 \text{Group}_{ij} + b_{1j} + \epsilon_{ij},$$

where  $b_{1j}$  is a random clinic effect. The random effect  $b_{1j}$  varies across clinics but, for a given clinic, is constant and shared by all patients belonging to that clinic. The inclusion of  $b_{1j}$  accounts for the clustering of patients within clinics, due perhaps to similarities in severity of illness and/or quality of care. The model explicitly accounts for the fact that some clinics have patients that respond higher (or lower) than patients in other clinics. However, the model assumes that the effect of treatment is the same across all clinics. This assumption can be relaxed by allowing the effect of treatment to vary among clinics

$$\begin{aligned} Y_{ij} &= \beta_1 + \beta_2 \text{Group}_{ij} + b_{1j} + b_{2j} \text{Group}_{ij} + \epsilon_{ij} \\ &= (\beta_1 + b_{1j}) + (\beta_2 + b_{2j}) \text{Group}_{ij} + \epsilon_{ij}. \end{aligned}$$

In this model the magnitude of the effect of treatment varies randomly across the different clinics. The average effect of treatment, when averaged over the population of clinics (and not simply those included in the trial), is  $\beta_2$ .

The example just presented involves randomizing patients (level 1) within clinics (level 2). In the language of experimental design, patients (level 1) are *nested* within clinics (level 2), but treatment is *crossed* with clinics because patients within each clinic are randomized to each treatment. Another very different type of design is one where patients (level 1) are nested within a clinic (level 2), but clinics are randomized to treatments, so that all patients from any given clinic receive the same treatment. In this case clinics are nested within treatment and not crossed. Formally, the same model just presented can be used for the analyses of these data, except that the effect of treatment can no longer vary randomly across the different clinics, since each clinic is assigned to only one treatment group. Note also that the treatment group variable,  $\text{Group}_{ij}$ , does not vary over  $i$  for fixed  $j$ , and hence can be replaced by  $\text{Group}_j$ . However, it is important to note that the nesting of clinics within treatment has a negative impact on efficiency of the treatment effect estimate, relative to a design with no nesting. This general principle will be illustrated later with analyses of the *Television, School and Family Smoking Prevention and Cessation Project*. In this study schools were randomized to treatments, and in the analysis both classroom and student variability were accounted for. The first design, where clinics are crossed with treatment, is generally more efficient than a design which does not stratify on clinic. The principle behind this is the same as that of a longitudinal study, where we can generally measure change more efficiently by using repeated measures on the same subject than by using a cross-sectional design.

So far our discussion of two-level models has very closely paralleled the description of the linear mixed effects model given in Chapter 8. In Chapter 8 we focused on models for two-level data where measurement occasions are the level 1 units and subjects are the level 2 units. However, it should be recognized that longitudinal data are simply a special case of two-level data and the linear mixed effects model given by (22.2) can be applied more broadly.

Finally, the two-level model given by (22.2) can also be written in terms of two models, one for each level of the hierarchy, using the two-stage formulation described in Section 8.4. That is, the two-level model can be expressed in terms of a level 1 model,

$$Y_{ij} = Z_{ij}\beta_j + \epsilon_{ij},$$

where  $\epsilon_{ij}$  are assumed to be independent across level 1 units, with mean zero and variance,  $\text{Var}(\epsilon_{ij}) = \sigma^2$ , and a level 2 model,

$$\beta_j = A_j\beta + b_j,$$

where  $b_j$  are assumed to vary independently across level 2 units, with mean zero and covariance,

$\text{Cov}(b_j) = G$ . Substituting the second model equation into the first yields (22.2)

$$\begin{aligned} Y_{ij} &= Z_{ij}(A_j\beta + b_j) + \epsilon_{ij} \\ &= (Z_{ij}A_j)\beta + Z_{ij}b_j + \epsilon_{ij} \\ &= X_{ij}\beta + Z_{ij}b_j + \epsilon_{ij}, \end{aligned}$$

where  $X_{ij} = Z_{ij}A_j$ . An advantage of specifying a multilevel model in terms of a series of models for each level of the hierarchy, rather than as a combined model, is that it becomes more transparent which covariates are operating at which level of the model. However, this does introduce some unnecessary restrictions on the model of the kind discussed in Section 8.4.

In summary, the two-level linear model given by (22.2) accounts for the clustering of the level 1 units by incorporating random effects at level 2. The model explicitly distinguishes two main sources of variation in the response: variation across level 2 units and variation across level 1 units (within level 2 units). The relative magnitude of these two sources of variability determines the degree of clustering in the data. The larger the variance of the level 2 random effects, relative to the level 1 (within level 2) variability, the greater is the degree of clustering. Next we describe how the linear mixed effects model can be generalized to three-level data structures; the extension to four or more levels follows directly.

## 22.3.2 Three-Level Linear Models

As mentioned earlier, there can be many levels in the data hierarchy. The extension of the two-level model given by (22.2) to three or more levels is very natural. With three-level data, there is clustering in the data that is assumed to be due to variation in the response across level 1, level 2, and level 3 units. Although the extension from 2 to 3 levels is conceptually straightforward, the description of the three-level model does require the introduction of additional notation that often obfuscates the salient features of the model. The basis of a three-level model is that variability in the response is accounted for by the introduction of random effects at all higher levels in the hierarchy (e.g., by allowing random variation in a subset of the regression parameters at both levels 2 and 3). The model explicitly distinguishes three sources of variation in the response: (1) variation across level 3 units, (2) variation across level 2 units (within level 3 units), and (3) variation across level 1 units (within level 2 units nested within level 3 units).

For a three-level data structure, let  $i$  index level 1 units,  $j$  index level 2 units, and  $k$  index level 3 units. We assume that there are  $n_3$  units at level 3 in the sample. Each of these clusters (for  $k = 1, \dots, n_3$ ) is composed of  $n_{2k}$  level 2 clusters, and each of these, in turn, is composed of  $n_{1jk}$  level 1 units. For example, consider a multi-center *longitudinal* clinical trial comparing two treatments (active drug versus placebo) conducted in 20 different centers or clinics. Patients in each clinic are measured at baseline and at three post-treatment occasions. In this example clinics are the level 3 units ( $k = 1, \dots, 20$ ), patients are the level 2 units ( $j = 1, \dots, n_{2k}$ ), and measurement occasions are the level 1 units ( $i = 1, \dots, 4$ ), where  $n_{2k}$  is the number of patients in the  $k^{\text{th}}$  clinic ( $n_3 = 20$  and  $n_{1jk} = 4$ , for all  $j$  and  $k$ ).

Let  $Y_{ijk}$  denote the response of the  $i^{\text{th}}$  level 1 unit within the  $j^{\text{th}}$  level 2 cluster within the  $k^{\text{th}}$  level 3 cluster. For example, in a multi-center longitudinal clinical trial,  $Y_{ijk}$  denotes the outcome at the  $i^{\text{th}}$  occasion for the  $j^{\text{th}}$  patient in the  $k^{\text{th}}$  clinic. Associated with each  $Y_{ijk}$  is a  $1 \times p$  (row) vector of covariates,  $X_{ijk}$ , with the covariates defined at different levels. A three-level model for  $Y_{ijk}$  is given by

$$(22.3) \quad Y_{ijk} = X_{ijk}\beta + Z_{ijk}^{(3)} b_k^{(3)} + Z_{ijk}^{(2)} b_{jk}^{(2)} + \epsilon_{ijk},$$

where  $Z_{ijk}^{(3)}$  is a design matrix for the random effects at level 3, formed from a subset of the appropriate columns of  $X_{ijk}$ , and  $Z_{ijk}^{(2)}$  is a design matrix for the random effects at level 2 (also formed from a subset of the appropriate columns of  $X_{ijk}$ ). In this notation the superscripts attached to  $b_k^{(3)}$  and  $b_{jk}^{(2)}$  denote the levels at which the random effects vary. In general, the design matrices for the random effects contain covariates that vary at lower levels than that of the corresponding random effects. That is,  $Z_{ijk}^{(3)}$ , the design matrix for  $b_k^{(3)}$ , will, in general, contain covariates that vary across level 2 and level 1 units.

To fix ideas, consider the example of a multi-center longitudinal clinical trial introduced earlier. A three-level model for the outcome is given by

$$Y_{ijk} = \beta_1 + \beta_2 t_{ijk} + \beta_3 (\text{Group}_{ij} \times t_{ijk}) + b_k^{(3)} + b_{jk}^{(2)} + \epsilon_{ijk},$$

where  $t_{ijk}$  denotes the time since baseline for the  $i^{\text{th}}$  observation on the  $j^{\text{th}}$  patient in the  $k^{\text{th}}$  clinic. In this model,  $b_k^{(3)}$  is a random clinic effect and  $b_{jk}^{(2)}$  is a random patient effect. The inclusion of the former accounts for the clustering of patients within clinics, while the inclusion of the latter accounts for the positive correlation among the repeated measures on the same patient. Additional random effects, at both levels 2 and 3, can easily be incorporated in the model (e.g., random slopes for time, and possibly random effects for the  $\text{Group}_{ij} \times t_{ijk}$  interaction).

The three-level model makes the following two assumptions about the different sources of variability:

1. The random effects  $b_k^{(3)}$  are assumed to be independent across level 3 units, with mean zero

and covariance,  $\text{Cov}(b_k^{(3)}) = G^{(3)}$ ; similarly the random effects  $b_{jk}^{(2)}$  are assumed to be independent across level 2 units, with mean zero and covariance,  $\text{Cov}(b_{jk}^{(2)}) = G^{(2)}$ . Random effects may be correlated within a given level, but not between levels.

2. The level 1 random components,  $\epsilon_{ijk}$ , are assumed to be independent across level 1 units, with mean zero and variance,  $\text{Var}(\epsilon_{ijk}) = \sigma^2$ . In addition the  $\epsilon_{ijk}$ 's are assumed to be independent of the random effects,  $b_k^{(3)}$  and  $b_{jk}^{(2)}$ .

That is, the random effects at the same level are, in general, correlated within units at that level but not between units; random effects at different levels are assumed to be independent of each other and of the level 1 random components,  $\epsilon_{ijk}$ . In principle, we can replace  $\epsilon_{ijk}$  in (22.3) with  $Z_{ijk}^{(1)} b_{ijk}^{(1)}$ , where  $b_{ijk}^{(1)}$  has mean zero and  $\text{Cov}(b_{ijk}^{(1)}) = G^{(1)}$ . This would allow for heterogeneity in the level 1 variability, with possible dependence of the level 1 variance on certain covariates. However, for the remainder of this discussion, we assume the simpler variance structure for the  $\epsilon_{ijk}$ , with  $\text{Var}(\epsilon_{ijk}) = \sigma^2$  (i.e., we assume  $Z_{ijk}^{(1)} = 1$  for all  $i, j$ , and  $k$ ).

In model (22.3) the regression parameters,  $\beta$ , are the fixed effects and describe the effects of covariates on the mean response (averaged over level 1, level 2, and level 3 units),

$$E(Y_{ijk}) = X_{ijk}\beta.$$

The three-level model given by (22.3) also describes the effects of covariates on the conditional mean response

$$E(Y_{ijk} | b_k^{(3)}) = X_{ijk}\beta + Z_{ijk}^{(3)} b_k^{(3)},$$

where the response is averaged over level 1 and level 2 units only, and the effects of covariates on the conditional mean response

$$E(Y_{ijk} | b_k^{(3)}, b_{jk}^{(2)}) = X_{ijk}\beta + Z_{ijk}^{(3)} b_k^{(3)} + Z_{ijk}^{(2)} b_{jk}^{(2)},$$

where the response is averaged over level 1 units only.

### 22.3.3 Estimation

The general specification of two- and three-level models that we have just described can be readily extended to more levels. The parameters of multilevel models are the fixed effects regression parameters,  $\beta$ , and the covariance (or variance) of the random effects at each level. Given estimates of the latter, predictions (empirical BLUPs) of the random effects at any level can also be obtained. For multilevel linear models, it is common to assume that the random components have multivariate normal distributions. For example, in the three-level model it is usually assumed that  $b_k^{(3)} \sim N(0, G^{(3)})$ ,  $b_{jk}^{(2)} \sim N(0, G^{(2)})$ , and  $\epsilon_{ijk} \sim N(0, \sigma^2)$ . Given these distributional assumptions, maximum likelihood (ML) estimation of the multilevel model parameters is relatively straightforward.

The ML estimate of  $\beta$  is obtained from the generalized least squares (GLS) estimator. For the two-level model, the GLS estimator has the same form as that given in Chapter 4 (albeit, with the indices  $i$  and  $j$  reversed). For the three-level model, the GLS estimator of  $\beta$  also has a closed-form expression and is given by

$$\hat{\beta} = \left\{ \sum_{k=1}^{n_3} (X_k' V_k^{-1} X_k) \right\}^{-1} \sum_{k=1}^{n_3} (X_k' V_k^{-1} Y_k),$$

where  $Y_k$  is a column vector, of length  $\sum_{j=1}^{n_{2k}} n_{1jk}$ , formed by stacking the responses for all second- and first-level units within the  $k^{th}$  cluster. Similarly,  $X_k$  is an  $(\sum_{j=1}^{n_{2k}} n_{1jk}) \times p$  matrix formed by stacking the covariates for all second- and first-level units within the  $k^{th}$  cluster. Finally,  $V_k$  is the covariance among observations on first- and second-level units within the  $k^{th}$  cluster and has a random effects covariance structure, expressed as a function of  $G^{(3)}$ ,  $G^{(2)}$ , and  $\sigma^2$  (and the corresponding design matrices for the random effects).

The restricted maximum likelihood (REML) estimates of  $G^{(3)}$ ,  $G^{(2)}$  and  $\sigma^2$  are obtained by maximizing the restricted log-likelihood with respect to  $G^{(3)}$ ,  $G^{(2)}$ , and  $\sigma^2$ . In general, it is not possible to write down simple, closed-form expressions for the REML estimators of  $G^{(3)}$ ,  $G^{(2)}$ , and  $\sigma^2$ ; instead, estimates must be obtained using iterative techniques. Once the REML estimates of  $G^{(3)}$ ,  $G^{(2)}$ , and  $\sigma^2$  have been obtained, the estimate of  $V_k(G^{(3)}, G^{(2)}, \sigma^2)$ , say  $V_k(\hat{g}^{(3)}, \hat{g}^{(2)}, \sigma^2)$ , is substituted into the generalized least squares estimator of  $\beta$  to obtain the REML estimate of  $\beta$ . REML estimation for multilevel linear models has been implemented in many major statistical software packages (e.g., PROC MIXED in SAS, the `lme` function in the `nlme` package in R and S-Plus, and the `xtmixed` command in Stata) and in stand-alone programs that have been specifically tailored for multilevel modeling (e.g., MLwiN and HLM).

## 22.3.4 Case Studies

Next we illustrate the main ideas by conducting analyses of two- and three-level data. The first example analyzes two-level data on fetal weight from a developmental toxicity study of laboratory mice exposed to ethylene glycol (EG). The data on the weights of live fetuses, nested within litters, are from an experiment conducted through the National Toxicology Program (NTP) (Price et al., 1985). The second example analyzes three-level data from a cluster-randomized trial to determine the efficacy of school-based interventions to prevent tobacco use. The data on seventh grade children, nested within classrooms, nested within schools are from the *Television, School, and Family Smoking Prevention and Cessation Project* (TVSFP) (Flay et al., 1995, Hedeker and Gibbons, 1994).

# Developmental Toxicity Study of Ethylene Glycol

Developmental toxicity studies of laboratory animals play a crucial role in the testing and regulation of chemicals and pharmaceutical compounds. Exposure to developmental toxicants typically causes a variety of adverse effects, such as fetal malformations and reduced fetal weight at term. In a typical developmental toxicity experiment, laboratory animals are assigned to increasing doses of a chemical or test substance. In this section we describe an analysis of data from a development toxicity study of ethylene glycol (EG). Ethylene glycol is a high-volume industrial chemical used in many applications. It is used as an antifreeze, as a solvent in the paint and plastics industries, and in the formulation of various types of inks. In a study of laboratory mice conducted through the National Toxicology Program (NTP), EG was administered at doses of 0, 750, 1500, or 3000 mg/kg/day to 94 pregnant mice (dams) over the period of major organogenesis, beginning just after implantation. (See Price et al., 1985, for additional details concerning the study design.) Following sacrifice, fetal weight and evidence of malformations were recorded for each live fetus. In our analysis of the data, we focus on the effects of dose on fetal weight; in Section 22.4 we present a complementary analysis that examines the effect of dose on fetal malformations.

Summary statistics (ignoring clustering in the data) for fetal weight for the 94 litters (composed of a total of 1028 live fetuses) are presented in [Table 22.1](#). Fetal weight decreases monotonically with increasing dose, with the average weight ranging from 0.972 (gm) in the control group to 0.704 (gm) in the group administered the highest dose. The decrease in fetal weight is not linear in increasing dose, but is approximately linear in increasing  $\sqrt{\text{dose}}$ .

**Table 22.1** Descriptive statistics on fetal weight from the ethylene glycol (EG) experiment.

Dose (mg/kg)	$\sqrt{\text{Dose}/750}$	Dams	Fetuses	Weight (gm)	
				Mean	St. Deviation <sup>a</sup>
0	0	25	297	0.972	0.098
750	1	24	276	0.877	0.104
1500	1.4	22	229	0.764	0.107
3000	2	23	226	0.704	0.124

<sup>a</sup>Calculated ignoring clustering.

The data on fetal weight from this experiment are clustered, with observations on the fetuses (level 1 units) nested within dams/litters (level 2 units). The litter sizes range from 1 to 16. Letting  $Y_{ij}$  denote the fetal weight of the  $i^{th}$  live fetus from the  $j^{th}$  litter, we considered the following model relating the fetal weight outcome to dose:

$$Y_{ij} = \beta_1 + \beta_2 d_j + b_j + \epsilon_{ij},$$

where  $d_j = \sqrt{\text{Dose}_j/750}$  is the square-root transformed dose administered to the  $j^{th}$  dam. The random effect  $b_j$  is assumed to vary independently across litters, with  $b_j \sim N(0, \sigma_2^2)$ . The errors,  $\epsilon_{ij}$ , are assumed to vary independently across fetuses (within a litter), with  $\epsilon_{ij} \sim N(0, \sigma_1^2)$ . Note that in a slight departure from the notation introduced previously, the first- and second-level variances are denoted by  $\sigma_1^2$  and  $\sigma_2^2$ , respectively. This model assumes that fetuses within a cluster are exchangeable and the positive correlation among the fetal weights is accounted for by their sharing a common random effect,  $b_j$ . The degree of clustering in the data can be expressed in terms of the intra-cluster (or intra-litter) correlation

$$\rho = \frac{\sigma_2^2}{\sigma_1^2 + \sigma_2^2}.$$

In [Table 22.2](#) the results of fitting the model to the fetal weight data are presented. The REML estimate of the regression parameter for (transformed) dose indicates that the mean fetal weight decreases with increasing dose. The estimated decrease in weight, comparing the highest dose group

to the control group, is 0.27 (or  $2 \times -0.134$ , with 95% confidence interval:  $-0.316$  to  $-0.220$ ). Of note, we calculated both model-based and empirical (or “sandwich”) standard errors and they were very similar, suggesting that the simple random effect structure for the clustering of fetal weights is adequate. The estimate of the intra-cluster correlation,  $\hat{\rho} = 0.57$ , indicates that there are moderate litter effects.

**Table 22.2** Fixed and random effects estimates for the fetal weight data from the ethylene glycol (EG) experiment.

Variable	Estimate	SE	Z
Intercept	0.984	0.016	61.32
$\sqrt{\text{Dose}/750}$	-0.134	0.012	-10.85
<b>Level 2 Variance:</b>			
$\sigma_2^2 (\times 100)$	0.726	0.119	
<b>Level 1 Variance:</b>			
$\sigma_1^2 (\times 100)$	0.556	0.026	

Finally, to assess the adequacy of the linear dose-response trend, we considered a model that included a quadratic effect of (transformed) dose. Both Wald and likelihood ratio tests of the quadratic effect of dose indicated that the linear trend is adequate for these data (Wald  $W^2 = 1.38$ , with 1 df,  $p > 0.20$ ; likelihood ratio  $G^2 = 1.37$ , with 1 df,  $p > 0.20$ ).

# Television School and Family Smoking Prevention and Cessation Project

Although smoking prevalence has declined among adults in recent decades, substantial numbers of young people begin to smoke and become addicted to tobacco. The *Television, School and Family Smoking Prevention and Cessation Project* (TVSFP) was a study designed to determine the efficacy of a school-based smoking prevention curriculum in conjunction with a television-based prevention program, in terms of preventing smoking onset and increasing smoking cessation (Flay et al., 1995). The study used a  $2 \times 2$  factorial design, with four intervention conditions determined by the cross-classification of a school-based social-resistance curriculum (CC: coded 1 = yes, 0 = no) with a television-based prevention program (TV: coded 1 = yes, 0 = no). Randomization to one of the four intervention conditions was at the school level, while much of the intervention was delivered at the classroom level. The TVSFP study is described in greater detail in Flay et al. (1995).

The original study involved 6695 students in 47 schools in Southern California. Our analysis focuses on a subset of 1600 seventh-grade students from 135 classes in 28 schools in Los Angeles. The response variable, a tobacco and health knowledge scale (THKS), was administered before and after randomization of schools to one of the four intervention conditions. The scale assessed a student's knowledge of tobacco and health.

We considered a linear model for the post-intervention THKS score, with the baseline or pre-intervention THKS score as a covariate. This model for the adjusted change in THKS scores included the main effects of CC and TV and the CC  $\times$  TV interaction. School and classroom effects were modeled by incorporating random effects at levels 3 and 2, respectively. Letting  $Y_{ijk}$  denote the post-intervention THKS score of the  $i^{th}$  student within the  $j^{th}$  classroom within the  $k^{th}$  school, our model is given by

$$Y_{ijk} = \beta_1 + \beta_2 \text{Pre-THKS} + \beta_3 \text{CC} + \beta_4 \text{TV} + \beta_5 \text{CC} \times \text{TV} + b_k^{(3)} + b_{jk}^{(2)} + \epsilon_{ijk},$$

where  $\epsilon_{ijk} \sim N(0, \sigma_1^2)$ ,  $b_{jk}^{(2)} \sim N(0, \sigma_2^2)$ , and  $b_k^{(3)} \sim N(0, \sigma_3^2)$ . Once again, in a slight departure from the notation introduced previously, the first-, second-, and third-level variances are denoted by  $\sigma_1^2$ ,  $\sigma_2^2$ , and  $\sigma_3^2$ , respectively.

The results of fitting this model to the data are presented in [Table 22.3](#). The REML estimates of the three sources of variability indicate that there is variability at both classroom and school levels, with almost twice as much variability among classrooms within a school as among schools. The correlation among the THKS scores for classmates (or children within the same classroom within the same school) is approximately 0.061 (or  $\frac{0.039+0.065}{0.039+0.065+1.602}$ ), while the correlation among the THKS scores for children from different classrooms within the same school is approximately 0.023 (or  $\frac{0.039}{0.039+0.065+1.602}$ ). The estimates of the fixed effects for the intervention conditions, when compared to their standard errors, indicate that neither the mass-media intervention (TV) nor its interaction with the social-resistance classroom curriculum (CC) has an impact on adjusted changes in the THKS scores from baseline. There is a significant effect of the social-resistance classroom curriculum, with children assigned to the social-resistance curriculum showing increased knowledge about tobacco and health. The estimate of the main effect of CC, in the model that excludes the CC  $\times$  TV interaction, is 0.47 (SE = 0.113,  $p < 0.0001$ ).

**Table 22.3** Fixed and random effects estimates for the THKS scores from the Television, School and Family Smoking Prevention and Cessation Project.

Variable	Estimate	SE	Z
Intercept	1.702	0.1254	13.57
Pre-Intervention THKS	0.305	0.0259	11.79
CC	0.641	0.1609	3.99
TV	0.182	0.1572	1.16
CC × TV	-0.331	0.2245	-1.47
<b>Level 3 Variance:</b>			
$\sigma_3^2$	0.039	0.0253	
<b>Level 2 Variance:</b>			
$\sigma_2^2$	0.065	0.0286	
<b>Level 1 Variance:</b>			
$\sigma_1^2$	1.602	0.0591	

The intra-cluster correlations at both the school and classroom levels are relatively small. The reader might be tempted to regard this as an indication that the clustering in these data is inconsequential. However, such a conclusion would be erroneous. Although the intra-cluster correlations are relatively small, they have a substantial impact on inferences concerning the effects of the intervention conditions. To illustrate this, consider the following model for the adjusted changes in THKS scores:

$$Y_{ijk} = \beta_1 + \beta_2 \text{Pre-THKS} + \beta_3 \text{CC} + \beta_4 \text{TV} + \beta_5 \text{CC} \times \text{TV} + \epsilon_{ijk},$$

where  $\epsilon_{ijk} \sim N(0, \sigma^2)$ . This model ignores clustering in the data at the classroom and school levels; it is a standard linear regression model and assumes independent observations and homogeneous variance. The results of fitting this model to the THKS scores are presented in [Table 22.4](#). The estimates of the fixed effects are similar to those reported in [Table 22.3](#). However, the model-based standard errors (assuming no clustering) are misleadingly small for the randomized intervention effects and lead to substantively different conclusions about the effects of the intervention conditions. This highlights an important lesson: the impact of clustering depends on both the magnitude of the intra-cluster correlation and the cluster size. For the data from the TVSFP, the cluster sizes vary from 1 to 13 classrooms within a school and from 2 to 28 students within a classroom. With relatively large cluster sizes, even very modest intra-cluster correlation can have a discernible impact on inferences.

**Table 22.4** Fixed effects estimates from analysis that ignores clustering in the THKS scores from the Television, School and Family Smoking Prevention and Cessation Project.

Variable	Estimate	SE	Z
Intercept	1.661	0.0844	19.69
Pre-Intervention THKS	0.325	0.0258	12.58
CC	0.641	0.0921	6.95
TV	0.199	0.0900	2.21
CC × TV	-0.322	0.1302	-2.47

## 22.4 MULTILEVEL GENERALIZED LINEAR MODELS

So far the discussion of multilevel models has focused on linear models for a continuous response, where clustering was accounted for through the introduction of random effects at different levels. Next we briefly describe how multilevel modeling can be extended to discrete response data. These models can be thought of as multilevel generalized linear models, and they extend in a natural way the conceptual approach described in Section 22.3. However, they differ in terms of assumptions concerning the distribution of observations at level 1. The level 1 observations are no longer required to have a normal distribution; instead, they are assumed to have a distribution belonging to the exponential family (e.g., Bernoulli or Poisson). We focus on models for two- and three-level data; the generalizations to more levels follow directly.

## 22.4.1 Two-Level Generalized Linear Models

The basic premise of multilevel generalized linear models is that clustering among units can be thought of as arising from their sharing a set of random effects. For example, with two-level binary data, it is assumed that the clustering of level 1 units (within level 2 units) can be accounted for by heterogeneity across level 2 clusters in a subset of the regression coefficients from a generalized linear model (e.g., a logistic regression model with randomly varying intercepts). Conditional on the random effects, the level 1 observations are assumed to be independent and with a distribution belonging to the exponential family (e.g., Bernoulli).

In our description of two-level generalized linear models we adopt the notation used earlier. Let  $Y_{ij}$  denote the response on the  $i^{th}$  level 1 unit in the  $j^{th}$  level 2 cluster; the response can be continuous, binary, or a count. Associated with each  $Y_{ij}$  is a  $1 \times p$  (row) vector of covariates,  $X_{ij}$ . We can formulate two-level models for discrete (and continuous) outcomes,  $Y_{ij}$ , using the familiar three-part specification of generalized linear mixed effects models outlined in Chapter 14:

1. We assume that the conditional distribution of each  $Y_{ij}$ , given a vector of random effects  $b_j$  (and the covariates), belongs to the exponential family of distributions and that  $\text{Var}(Y_{ij}|b_j) = v\{E(Y_{ij}|b_j)\} \phi$ , where  $v(\cdot)$  is a known variance function, a function of the conditional mean,  $E(Y_{ij}|b_j)$ , and  $\phi$  is a scale or dispersion parameter. In addition, given the random effects  $b_j$ , it is assumed that the  $Y_{ij}$  are independent of one another.
2. The conditional mean of  $Y_{ij}$  is assumed to depend on fixed and random effects via the following linear predictor:

$$\eta_{ij} = X_{ij}\beta + Z_{ij}b_j,$$

with

$$g\{E(Y_{ij}|b_j)\} = \eta_{ij} = X_{ij}\beta + Z_{ij}b_j$$

for some known link function,  $g(\cdot)$ .

3. The random effects are assumed to have some probability distribution. In principle, any multivariate distribution can be assumed for the  $b_j$ ; in practice, for computational convenience, the random effects are usually assumed to have a multivariate normal distribution, with zero mean and covariance matrix,  $G$ .

These three components completely specify a broad class of two-level generalized linear models for different types of responses. Next, to clarify the main ideas, we consider two examples of multilevel generalized linear models in greater detail.

# Example 1: Two-Level Generalized Linear Model for Counts

Consider a study comparing cross-national rates of skin cancer and the factors (e.g., climate, economic and social factors, regional differences in diagnostic procedures) that influence variability in the rates of disease. Suppose that we have counts of the number of cases of skin cancer in a set of well-defined regions, indexed by  $i$ , within countries, indexed by  $j$ . Let  $Y_{ij}$  be a count of the number of individuals who develop skin cancer within the  $i^{th}$  region of the  $j^{th}$  country during a given period of time (e.g., 5 years). The resulting counts have a two-level structure with regional units at the lower level (level 1 units) nested within countries (level 2 units). Usually the analysis of count data requires knowledge of the denominator, the population at risk. That is, the *rate* at which the disease occurs is of more direct interest than the corresponding count.

Counts are often modeled as Poisson random variables using a log link function. This motivates the following illustration of a two-level generalized linear model for  $Y_{ij}$  given by the three-part specification:

1. Conditional on a vector of random effects  $b_j$ , the  $Y_{ij}$  are assumed to be independent observations from a Poisson distribution, with  $\text{Var}(Y_{ij}|b_j) = E(Y_{ij}|b_j)$ , (i.e.,  $\phi = 1$ ).
2. The conditional mean of  $Y_{ij}$  depends on fixed and random effects via the following log link function,

$$\log \{E(Y_{ij}|b_j)\} = \log(T_{ij}) + X_{ij}\beta + Z_{ij}b_j,$$

where  $T_{ij}$  is the population at risk in the  $i^{th}$  region of the  $j^{th}$  country and  $\log(T_{ij})$  is an *offset*.

3. The random effects are assumed to have a multivariate normal distribution, with zero mean and covariance matrix  $G$ .

This is an example of a two-level log-linear model that assumes a linear relationship between the log rate of disease occurrence and the covariates.

## Example 2: Two-Level Generalized Linear Model for Binary Responses

Consider a study of men with newly diagnosed prostate cancer. The study is designed to evaluate the factors that determine physician recommendations for surgery (radical prostatectomy) versus radiation therapy. In particular, it is of interest to determine the relative importance of patient factors (e.g., patient's age, level of prostate specific antigen) and physician factors (e.g., specialty training, years of experience) on physician recommendations for treatment. Many patients in the study seek the recommendation of the same physician. As a result patients (level 1 units) are nested within physicians (level 2 units). For each patient we have a binary outcome denoting the physician recommendation (surgery versus radiation therapy).

Let  $Y_{ij}$  be the binary response, taking values 0 and 1 (e.g., denoting surgery or radiation therapy) for the  $i^{th}$  patient of the  $j^{th}$  physician. An illustrative example of a two-level logistic model for  $Y_{ij}$  is given by the following three-part specification:

1. Conditional on a single random effect  $b_j$ , the  $Y_{ij}$  are independent and have a Bernoulli distribution, with  $\text{Var}(Y_{ij}|b_j) = E(Y_{ij}|b_j)\{1 - E(Y_{ij}|b_j)\}$ , (i.e.,  $\phi = 1$ ).
2. The conditional mean of  $Y_{ij}$  depends on fixed and random effects via the following linear predictor:

$$\eta_{ij} = X_{ij}\beta + b_j,$$

with

$$\log \left\{ \frac{\Pr(Y_{ij} = 1|b_j)}{\Pr(Y_{ij} = 0|b_j)} \right\} = \eta_{ij} = X_{ij}\beta + b_j.$$

That is, the conditional mean of  $Y_{ij}$  is related to the linear predictor by a logit link function.

3. The single random effect  $b_j$  is assumed to have a univariate normal distribution, with zero mean and variance  $g_{11}$ .

In this example the model is a simple two-level logistic regression model with randomly varying intercepts. In principle, the linear predictor can include additional random effects. However, some caution must be exercised because there is usually not much information in binary data to estimate more than a single variance component unless the number of level 1 units is relatively large.

In Section 11.3 we discussed how the logistic regression model for binary data can also be developed through the notion of an underlying latent variable distribution. The same notion can be applied to multilevel models for binary responses. Suppose that  $L_{ij}$  is a latent (i.e., unobserved) continuous variable and that a positive response is observed only when  $L_{ij}$  exceeds some threshold. Consider the following two-level linear model for  $L_{ij}$ ,

$$L_{ij} = X_{ij}\beta + Z_{ij}b_j + \epsilon_{ij},$$

where the random effects are assumed to have a multivariate normal distribution, with mean zero and covariance matrix,  $G$ , and the  $\epsilon_{ij}$  are assumed to have a standard logistic distribution, with mean zero and variance  $\pi^2/3$ . Without loss of generality, we can assume the threshold for categorizing  $L_{ij}$  is zero, with

$$Y_{ij} = 1 \text{ if } L_{ij} > 0,$$

$$Y_{ij} = 0 \text{ if } L_{ij} \leq 0.$$

Then the relationship between  $Y_{ij}$  and  $L_{ij}$  results in a logistic regression model for  $\Pr(Y_{ij} = 1|b_j)$ . That is, the two-level linear model for  $L_{ij}$  with standard logistic errors,

$$L_{ij} = X_{ij}\beta + Z_{ij}b_j + \epsilon_{ij},$$

implies the two-level logistic regression model for  $Y_{ij}$ ,

$$\log \left\{ \frac{\Pr(Y_{ij} = 1|b_j)}{\Pr(Y_{ij} = 0|b_j)} \right\} = X_{ij}\beta + Z_{ij}b_j.$$

Using the notion of an underlying latent variable distribution, we can then compare the magnitudes of the between-cluster and within-cluster sources of variability of the  $L_{ij}$ . For example, in a two-level logistic regression model with a single random effect  $b_j$  (with variance  $g_{11}$ ), the relative magnitudes of the between-cluster and within-cluster sources of variability can be summarized in terms of the intra-cluster correlation

$$\rho = \text{Corr}(L_{ij}, L_{i'j}) = \frac{g_{11}}{g_{11} + \pi^2/3}; \quad (\text{for } i \neq i').$$

Note that  $\rho$  is the marginal correlation among the latent variables,  $L_{ij}$ ; it is not the marginal correlation among the binary variables,  $Y_{ij}$ .

Although in both of the examples of two-level generalized linear models we have chosen canonical link functions to relate the conditional mean of  $Y_{ij}$ , to  $\eta_{ij}$ , in principle, any suitable link function can be selected. These two examples are intended to be purely illustrative. They demonstrate how the choices of the three components might differ according to the type of response variable.

So far our discussion of two-level models has closely paralleled the description of the generalized linear mixed effects model given in Chapter 14 (albeit with the indices  $i$  and  $j$  reversed). In Chapter 14 we focused on two-level models where measurement occasions are the level 1 units and subjects are the level 2 units; this is a special case of two-level data. However, the methods in Chapter 14 can be applied more broadly to different types of two-level data and also extend naturally to more than two levels.

## 22.4.2 Three-Level Generalized Linear Models

The extension of two-level generalized linear models to three or more levels is straightforward and follows from the previous sections. With three-level data the variability of the response is accounted for by the introduction of random effects at both levels 2 and 3. For a three-level data structure, we adopt the same notation as in Section 22.3, except that the response can be continuous, binary, or a count. Let  $Y_{ijk}$  denote the response on the  $i^{th}$  level 1 unit within the  $j^{th}$  level 2 cluster within the  $k^{th}$  level 3 cluster. Associated with each  $Y_{ijk}$  is a  $1 \times p$  (row) vector of covariates,  $X_{ijk}$ , with the covariates defined at different levels. A three-level generalized linear model for  $Y_{ijk}$  is given by the following three-part specification:

1. We assume that the conditional distribution of each  $Y_{ijk}$ , given vectors of random effects,  $b_k^{(3)}$  and  $b_{jk}^{(2)}$  (defined at levels 3 and 2, respectively), belongs to the exponential family of distributions and that  $\text{Var}(Y_{ijk}|b_k^{(3)}, b_{jk}^{(2)}) = v\{E(Y_{ijk}|b_k^{(3)}, b_{jk}^{(2)})\} \phi$ , where  $v(\cdot)$  is a known variance function, a function of the conditional mean,  $\{E(Y_{ijk}|b_k^{(3)}, b_{jk}^{(2)})\}$ , and  $\phi$  is a scale or dispersion parameter. In addition, given  $b_k^{(3)}$  and  $b_{jk}^{(2)}$ , it is assumed that the  $Y_{ijk}$  are independent of one another.
2. The conditional mean of  $Y_{ijk}$  is assumed to depend on fixed and random effects via the following linear predictor:

$$\eta_{ijk} = X_{ijk}\beta + Z_{ijk}^{(3)} b_k^{(3)} + Z_{ijk}^{(2)} b_{jk}^{(2)},$$

with

$$g\{E(Y_{ijk}|b_k^{(3)}, b_{jk}^{(2)})\} = \eta_{ijk} = X_{ijk}\beta + Z_{ijk}^{(3)} b_k^{(3)} + Z_{ijk}^{(2)} b_{jk}^{(2)},$$

for some known link function,  $g(\cdot)$ .

3. The random effects are assumed to have multivariate normal distributions, with mean zero and covariance matrices,  $\text{Cov}(b_k^{(3)}) = G^{(3)}$  and  $\text{Cov}(b_{jk}^{(2)}) = G^{(2)}$ . Although random effects may be correlated within a given level, random effects at different levels are assumed to be independent of each other.

These three components completely specify a broad class of three-level generalized linear models.

## 22.4.3 Estimation

The multilevel generalized linear models described in the previous section fully specify the joint distribution of the responses at level 1 and the random effects at all higher levels. As a result we can base estimation and inference on the likelihood function. However, unlike the case with a continuous response assumed to have a normal distribution, maximum likelihood (ML) estimation for multilevel generalized linear models is not straightforward and will, in general, require numerical quadrature.

For example, for three-level data, inference about  $\beta$ ,  $G^{(2)}$  and  $G^{(3)}$  is based on the marginal likelihood function. The marginal likelihood can be expressed as the product of the probability density functions,  $f(y_k)$ , for the level 3 units. Specifically, the marginal log-likelihood function is given by the following sum:

$$\sum_{k=1}^{n_3} \log f(y_k),$$

where  $f(y_k)$  can be obtained by recognizing that observations on level 2 units (within level 3 units) are conditionally independent of one another given the level 3 random effects,  $b^{(3)}$ ; similarly, observations on level 1 units (within level 2 units) are conditionally independent of one another given the level 3 and level 2 random effects,  $b^{(3)}$  and  $b^{(2)}$ . The  $k^{\text{th}}$  level 3 unit's contribution to the likelihood function is

$$f(y_k) = \int \prod_{j=1}^{n_{2k}} \left\{ \int \prod_{i=1}^{n_{1jk}} f(y_{ijk} | b_k^{(3)}, b_{jk}^{(2)}) f(b_{jk}^{(2)}) db_{jk}^{(2)} \right\} f(b_k^{(3)}) db_k^{(3)},$$

where  $f(b_{jk}^{(2)})$  and  $f(b_k^{(3)})$  denote the multivariate normal distributions for the random effects at levels 2 and 3, respectively. The ML estimates of  $\beta$ ,  $G^{(2)}$  and  $G^{(3)}$  are simply those values of  $\beta$ ,  $G^{(2)}$  and  $G^{(3)}$  that maximize the marginal log-likelihood function.

The primary reason for displaying the expression given above is to highlight that multivariate integrals must be evaluated to compute the marginal log-likelihood. That is, the log-likelihood function is obtained by integrating out or averaging over the distributions of the random effects,  $b_{jk}^{(2)}$  and  $b_k^{(3)}$ . Because integrals (denoting averaging over the distribution of the random effects) appear in the log-likelihood function, there are no simple, closed-form solutions. Instead, numerical integration techniques, for instance, Gaussian quadrature, are required for maximizing the log-likelihood function. ML estimation, using Gaussian quadrature, for two-level and higher-level generalized linear models is implemented in some of the major statistical software packages (e.g., PROC GLIMMIX in SAS, the `glmer` function in the `lme4` package in R, and the `xtmelogit` and `xtmepoisson` commands in Stata). Various alternative approximations to ML estimation for the extensions to three or more levels are implemented in more specialized, stand-alone programs that have been specifically developed for multilevel modeling (e.g., MLwiN and HLM).

## 22.4.4 Case Studies

We illustrate the main ideas underlying multilevel modeling by conducting analyses of two-level data where the observations at level 1 are counts and binary outcomes. The first example analyzes two-level count data from a study of malignant melanoma mortality and ultraviolet (UV) radiation exposure. The second example analyzes two-level data on fetal malformations, a binary outcome, from the developmental toxicity study of ethylene glycol (EG) described in Section 22.3. For the latter, we present a traditional multilevel analysis of the fetal malformation data and contrast the results with those obtained from a marginal model that accounts for clustering in the data in a different way.

# Malignant Melanoma Mortality and Ultraviolet Light Exposure

In a study of the effects of ultraviolet (UV) light exposure on malignant melanoma mortality (Langford et al., 1998), counts of the number of deaths due to malignant melanoma were recorded for males of all ages in the United Kingdom. The counts of the number of deaths between 1975 and 1980 were aggregated over areas that correspond to counties or shires; hereafter referred to as counties. Data were collected on 70 counties nested within 11 regions of the United Kingdom. The resulting data structure is multilevel, with counties at level 1 (indexed by  $i$ ) nested within regions at level 2 (indexed by  $j$ ). The main predictor of interest is exposure to ultraviolet light in the B band (UVB). An index of UVB dose reaching the earth's surface was calculated for each county. The mean UVB index in the United Kingdom was 10.9, with a standard deviation of 1.5.

Let  $Y_{ij}$  denote the count of the number of deaths in the  $i^{th}$  county in the  $j^{th}$  region due to malignant melanoma. Within a given region, we assume  $Y_{ij}$  has a Poisson distribution to account for level 1 variation in the counts. Variation in the counts across regions is accounted for by the inclusion of a random region effect,  $b_j$ . That is, conditional on a random region effect  $b_j$ , the counts are assumed to have a Poisson distribution with conditional mean related to UVB dose via a log link function,

$$\log\{E(Y_{ij}|b_j)\} = \log(T_{ij}) + \beta_1 + \beta_2 \text{UVB}_{ij} + b_j,$$

where  $\text{UVB}_{ij}$  is the UVB index (centered at the mean UVB index in the United Kingdom) in the  $i^{th}$  county of the  $j^{th}$  region. For each county,  $T_{ij}$  is the number of deaths that would be expected were U.K. national age- and gender-specific death rates to apply to the population of the county. Note that  $T_{ij}$  is known and  $\log(T_{ij})$  is an *offset* in this model. The ratio of the observed number of deaths to the expected number of deaths,  $Y/T$ , is referred to as the *standardized mortality ratio* (SMR) in each county. Our model assumes a linear relationship between the log SMR due to malignant melanoma and county-level UV radiation exposure. Finally, we assume the random region effect has a normal distribution,  $b_j \sim N(0, \sigma^2)$ .

The results of fitting this model to the UK malignant melanoma mortality data, using maximum likelihood estimation, are presented in [Table 22.5](#). There is a significant positive relationship between the SMRs and exposure to UVB. Recall that the standard deviation of UVB in the United Kingdom is 1.5. Therefore the estimated effect of UVB dose indicates that the SMR is approximately 1.5 times larger (or  $e^{3 \times 0.13}$ , with 95% confidence interval: 1.25 to 1.74) when comparing a county with UVB index 1 standard deviation above the UK average to a county with UVB index 1 standard deviation below. Finally, in interpreting these results, it should be remembered that the UVB covariate used here is simply an index of exposure for each county; it is the potential, not the actual, UVB dose experienced by the population of residents in each county.

**Table 22.5** Fixed and random effects estimates for the malignant melanoma mortality data for males in the United Kingdom.

Variable	Estimate	SE	Z
Intercept	-0.0365	0.0352	-1.04
UVB	0.1301	0.0279	4.67
<b>Level 2 Variance:</b>			
$\sigma^2 \times 100$	0.6222	0.5087	

*Note:* ML estimation is based on 50-point adaptive Gaussian quadrature.

# Developmental Toxicity Study of Ethylene Glycol

Next we consider a two-level logistic regression model for binary data on fetal malformations from the developmental toxicity study of ethylene glycol (EG). Recall that in this study, EG was administered (0, 750, 1500, or 3000 mg/kg/day) to 94 pregnant mice (dams) over the period of major organogenesis, beginning just after implantation (Price et al., 1985). Following sacrifice, each live fetus was examined for evidence of malformations, recorded as present or absent. The primary question of scientific interest is the effect of dose on fetal malformations.

Summary statistics (ignoring clustering in the data) for fetal malformations for the 94 litters (composed of a total of 1028 live fetuses) are presented in [Table 22.6](#). The percentage of fetal malformations increases monotonically with increasing dose, with less than 1 % in the control group and almost 60% in the group administered the highest dose (3000 mg/kg/day).

**Table 22.6** Descriptive statistics on fetal malformations from the ethylene glycol (EG) experiment.

Dose (mg/kg/day)	Dams	Fetuses	Fetal Malformations	
			Number	Percentage
0	25	297	1	0.34
750	24	276	26	9.42
1500	22	229	89	38.86
3000	23	226	129	57.08

Letting  $Y_{ij} = 1$  denote the presence of fetal malformations in the  $i^{th}$  live fetus from the  $j^{th}$  litter (and  $Y_{ij} = 0$  otherwise), we considered the following logistic model relating the log odds of fetal malformations to a linear effect of dose:

$$\text{logit}\{E(Y_{ij}|b_j)\} = \beta_1 + \beta_2 d_j + b_j,$$

where  $d_j = \text{Dose}_j/750$  denotes the dose (in units of 750 mg/kg/day) administered to the  $j^{th}$  dam (cluster). The random effect  $b_j$  is assumed to vary independently across litters, with  $b_j \sim N(0, \sigma_b^2)$ . This model assumes that fetuses within a litter are exchangeable and the positive association among the fetal malformation outcomes is accounted for by their sharing a common random effect,  $b_j$ .

The results of fitting the model to the fetal malformation data, using maximum likelihood estimation, are presented in [Table 22.7](#). The estimated regression parameter for dose indicates that the log odds of malformation increases with increasing dose. In particular, the odds ratio for malformation, comparing the highest dose group to the control group, is 209.2 (or  $e^{4 \times 1.336}$ , with 95% confidence interval: 56.1 to 779.9). This provides overwhelming evidence of the increased risk of malformations at the highest dose of EG. The odds ratio for malformations, comparing the lowest dose group to the control group, is 3.80 (or  $e^{1.336}$ , with 95% confidence interval: 2.75 to 5.26). Finally, the estimate of  $\sigma_b^2$  indicates that there are moderate litter effects, with heterogeneity across dams in the underlying risk of producing fetuses with malformations. For example, in the control group, 95% of dams have a risk of producing fetuses with malformations between 0% and 22%. Alternatively, if we appeal to the notion of a latent variable distribution and assume an underlying two-level linear model for the latent variable with standard logistic errors, the estimated intra-cluster correlation is

$$\hat{\rho} = \frac{\hat{\sigma}_b^2}{\hat{\sigma}_b^2 + \pi^2/3} = \frac{2.517}{2.517 + 3.290} = 0.43.$$

**Table 22.7** Fixed and random effects estimates for the fetal malformation data from the ethylene glycol (EG) experiment.

Variable	Estimate	SE	Z
Intercept	-4.360	0.440	-9.92
Dose / 750	1.336	0.166	8.06
<b>Level 2 Variance:</b>			
$\sigma_b^2$	2.517	0.685	

*Note:* ML estimation is based on 50-point adaptive Gaussian quadrature.

For illustrative purposes we also consider a marginal logistic regression model relating the log odds of fetal malformations to a linear effect of dose

$$\text{logit}\{E(Y_{ij})\} = \beta_1 + \beta_2 d_j,$$

and account for the intra-litter association by a common log odds ratio,

$$\log \{\text{OR}(Y_{ij}, Y_{ik})\} = \log \left\{ \frac{\Pr(Y_{ij} = 1, Y_{ik} = 1) \Pr(Y_{ij} = 0, Y_{ik} = 0)}{\Pr(Y_{ij} = 1, Y_{ik} = 0) \Pr(Y_{ij} = 0, Y_{ik} = 1)} \right\}.$$

This can be thought of as a marginal model analogue of the random intercepts model for the within-litter association. The results of fitting this marginal model, using GEE methods, are presented in [Table 22.8](#). The estimate of the effect of dose indicates that the odds ratio for malformations, comparing the highest dose group to the control group, is 46.4 (or  $e^{4 \times 0.960}$ , with 95% confidence interval: 21.3 to 101.2). These results also provide strong evidence of the increased risk of malformations at the highest dose of EG. The within-litter odds ratio of 4.26 (or  $e^{1.447}$ ) indicates that there is clustering in the fetal malformations data.

**Table 22.8** Estimates of regression parameters from marginal model for the fetal malformation data from the ethylene glycol (EG) experiment.

Variable	Estimate	SE	Z
Intercept	-3.190	0.220	-14.53
Dose / 750	0.960	0.099	9.66
Log Odds Ratio	1.447	0.221	6.56

Note that the estimated effect of dose in the marginal model is discernibly smaller than that reported in [Table 22.7](#). This should not be too surprising given the important distinctions between the regression parameters in marginal models and generalized linear mixed effects models that were highlighted in Chapter 16. The regression parameters for dose in the two models have quite different interpretations. In the logistic regression model with a random litter effect,  $\beta_2$  describes the change in the risk of producing fetuses with malformations for *any given dam*; this change in the risk for a single-unit change in dose depends on  $b_j$ , a specific dam's random effect or underlying propensity for producing fetuses with malformations. In the marginal model,  $\beta_2$  describes changes in the prevalence of fetal malformations when sub-populations of dams exposed to different doses of EG are compared. Although both models account for clustering in the data, the targets of inference are different and the two analyses address distinct scientific questions.

## 22.5 SUMMARY

In previous sections we described models for data with a hierarchical structure, where lower-level units are nested within higher-level units. The dominant approach for modeling such data is regression models where random effects are introduced at different levels. A central feature of multilevel modeling is the incorporation of covariates that can be measured at any level of the hierarchy, thereby allowing the effects at each level to be disentangled. By combining covariates that have been measured at different levels within a single regression model, their relative importance can be determined. For example, multilevel models can address questions about the effects of individuals' characteristics (e.g., disease severity) while adjusting for their context (e.g., being treated at a large university teaching hospital versus a rural clinic).

Multilevel data can be challenging to analyze for at least two main reasons. First, the covariates can be measured at different levels, and the same covariate can operate at many different levels. As a result somewhat greater care is required in the interpretation of regression parameters in multilevel models. It is not always transparent how best to combine covariates measured at different levels within a single model so that the regression parameters have useful interpretations. In our brief overview of multilevel models, we have not touched on this important topic; for more information, the interested reader is directed to the references at the end of this chapter.

The second challenge in the analysis of multilevel data is how best to account for the clustering that can arise at different levels of the hierarchy. In the multilevel modeling literature the dominant approach for accounting for the intra-cluster correlations at different levels is via the introduction of random effects at different levels. This gives rise to mixed effects models that can be extended in a very natural way to any number of levels of clustering in the data. For linear models, this is certainly a very natural way to account for clustering. However, for generalized linear models for discrete data, it does raise subtle issues concerning the interpretation of the fixed effect regression parameters and questions about what is the relevant target of inference. The same issues that were given an airing in Chapter 16 apply equally to multilevel models for discrete data. These issues were highlighted in our analyses of the two-level clustered data on fetal malformations in Section 22.4, where the estimated effect of dose was discernibly different depending on how the clustering was accounted for. The estimates of the effect of dose from the marginal and random effects logistic regression models differed because the corresponding regression parameters have distinct interpretations and address somewhat different scientific questions. In general, the fixed effects parameters in a two-level model for discrete data represent changes in the (transformed) mean response, for a single-unit change in the corresponding covariate, *for any given level 2 unit*. In Chapters 14 and 16, in the context of longitudinal data, these regression coefficients were referred to as "subject-specific"; here they are "cluster-specific" and describe covariate effects for an individual cluster. In contrast, the regression parameters in a marginal model represent changes in the (transformed) mean response when sub-populations defined by different values of the corresponding covariate are compared. The regression parameters in marginal models address the dependence of the population-averaged response (where averaging is over all possible units in the hierarchy) on the covariates. These regression parameters do not have any direct interpretation for an individual cluster when there is heterogeneity across clusters.

Although much of the multilevel literature on the analysis of discrete data is dominated by the use of generalized linear mixed effects models, we note that marginal models can also be used to account for clustering at different levels. All the issues discussed in Chapter 16 for two-level longitudinal data apply equally to two-level and higher-level data more broadly defined. In general, the choice between the two classes of models should not be driven by the availability of software for multilevel modeling but on the basis of careful thought about the questions of scientific interest.

## 22.6 FURTHER READING

There is an extensive literature on multilevel models that appears in the statistical, psychometric, and educational literature. A comprehensive description of multilevel models, and their application to a wide range of problems, can be found in the books by Raudenbush and Bryk (2002), Longford (1993), and Goldstein (2003). For readers who find the level of mathematical difficulty in these books too challenging, the books by Hox (2002), Kreft and De Leeuw (1998), and Snijders and Bosker (1999) provide a more introductory and accessible presentation of similar topics targeted at empirical researchers. An engaging and non-technical introduction to multilevel modeling can be found in the excellent text by Gelman and Hill (2007).

For illustrations of the application of multilevel models in the biomedical and health sciences, we recommend the edited volume of articles in Leyland and Goldstein (2001) and the review articles by Sullivan et al. (1999), Goldstein et al. (2002), and Subramanian et al. (2003).

## Bibliographic Notes

Although most of the statistical literature on marginal models has focused on two-level data, Qaqish and Liang (1992) discuss the use of marginal models for multilevel binary data, with multiple levels of nesting.

Daniels and Gatsonis (1999) describe multilevel modeling in a Bayesian framework; also see Lindley and Smith (1972), Zeger and Karim (1991), Browne et al. (2002), Carlin and Louis (2000), and Chapters 15 and 16 of Gelman et al. (2003). Bayesian methods for multilevel modeling can be implemented using the publicly available software WinBUGS (Spiegelhalter et al., 1999).

## *Appendix A*

### ***Gentle Introduction to Vectors and Matrices***

We present a very brief introduction to vectors and matrices, intended for readers with no prior exposure to matrix algebra. Specifically, we cover basic definitions and summarize some of the main properties of vectors and matrices. Vectors and matrices allow us to perform common mathematical operations (e.g., addition, subtraction, and multiplication) on a collection of numbers; they also facilitate the description of statistical methods for multivariate data. Our primary motivation for using them is the conciseness and compactness with which statistical techniques for analyzing longitudinal data can be presented when expressed in terms of vectors and matrices.

Mastery of the material presented in this section is a prerequisite for understanding the statistical methods for longitudinal data described in the book. Although we do not assume a profound understanding of matrix algebra, vectors and matrices are used extensively throughout the book to simplify notation and the reader is required to have some basic facility with the addition and multiplication of vectors and matrices.

# Basic Concepts and Definitions

A *matrix* is a rectangular array of elements (e.g., numbers), arranged in rows and columns. For example,

$$A = \begin{pmatrix} 2 & 7 & 11 & 5 \\ 4 & 9 & 3 & 1 \\ 13 & 8 & 2 & 6 \end{pmatrix}$$

is a matrix with three rows and four columns. The *element* or *entry* in the  $i^{th}$  row and the  $j^{th}$  column of the matrix is referred to as the  $(i, j)^{th}$  element of the matrix. For example, the entry in the  $2^{nd}$  row and  $3^{rd}$  column of  $A$  is 3. If we let  $a_{ij}$  denote the element in the  $i^{th}$  row and the  $j^{th}$  column of the matrix  $A$ , then

$$\begin{aligned} a_{11} &= 2, & a_{12} &= 7, & a_{13} &= 11, & a_{14} &= 5; \\ a_{21} &= 4, & a_{22} &= 9, & a_{23} &= 3, & a_{24} &= 1; \\ a_{31} &= 13, & a_{32} &= 8, & a_{33} &= 2, & a_{34} &= 6. \end{aligned}$$

The subscripts on the element  $a_{ij}$  denote its position in the  $i^{th}$  row and the  $j^{th}$  column of the matrix  $A$ .

The *dimension* of a matrix is the number of rows and columns in the matrix. By convention, the number of rows is listed first, and then the number of columns. Thus we refer to the matrix  $A$  above as being a  $3 \times 4$ , or a “3 by 4”, matrix.

A *vector* is a special kind of matrix, having either a single row or a single column. For example,

$$V = \begin{pmatrix} 2 \\ 4 \\ 9 \\ 7 \\ 3 \end{pmatrix}$$

is a  $5 \times 1$  (column) vector. Since the dimension of a vector corresponds to the number of elements in the vector, the dimension of a vector is often loosely referred to as its *length*.<sup>1</sup>

Finally, a *scalar* is a single element (e.g., a single number), and hence can be treated either as a single-element vector or as a  $1 \times 1$  matrix.

# Transpose

The *transpose* is a function that interchanges the rows and columns of a matrix. That is, the first row becomes the first column, the second row becomes the second column, and so on. By convention, the transpose of a matrix  $A$  is denoted  $A'$  (or “ $A$  prime”). (Note that in some texts, a superscript  $T$ , instead of a prime, is used to denote the transpose of a matrix, e.g.,  $A^T$ .)

For example, consider the  $3 \times 4$  matrix  $A$ ,

$$A = \begin{pmatrix} 2 & 7 & 11 & 5 \\ 4 & 9 & 3 & 1 \\ 13 & 8 & 2 & 6 \end{pmatrix}.$$

The transpose of  $A$ ,

$$A' = \begin{pmatrix} 2 & 4 & 13 \\ 7 & 9 & 8 \\ 11 & 3 & 2 \\ 5 & 1 & 6 \end{pmatrix},$$

is the  $4 \times 3$  matrix with rows and columns interchanged. Similarly, since a vector is a matrix with either a single row or column, if

$$V = \begin{pmatrix} 2 \\ 4 \\ 9 \\ 7 \\ 3 \end{pmatrix}, \quad \text{then } V' = (2 \ 4 \ 9 \ 7 \ 3).$$

Examples of vectors and matrices that play key roles in the analysis of longitudinal data are the *response vector*, often denoted  $Y$ , and the *covariate matrix*, often denoted  $X$ . For example, consider the following data from a subject participating in a longitudinal clinical trial. In this trial, the subject was assigned to the placebo group (Group = 0 if assigned to placebo, Group = 1 if assigned to active treatment) and four repeated measures of blood lead levels were obtained at baseline (or week 0), week 1, week 4, and week 6:

**Blood Lead Treatment Group Week**

30.8	0	0
26.9	0	1
25.8	0	4
23.8	0	6

If we let  $Y$  denote the vector of repeated measurements of the response variable, then

$$Y = \begin{pmatrix} Y_1 \\ Y_2 \\ Y_3 \\ Y_4 \end{pmatrix} = \begin{pmatrix} 30.8 \\ 26.9 \\ 25.8 \\ 23.8 \end{pmatrix}.$$

Similarly we can let  $X$  denote a matrix of covariates associated with the vector of repeated measurements, with

$$X = \begin{pmatrix} X_{11} & X_{12} & X_{13} \\ X_{21} & X_{22} & X_{23} \\ X_{31} & X_{32} & X_{33} \\ X_{41} & X_{42} & X_{43} \end{pmatrix} = \begin{pmatrix} 1 & 0 & 0 \\ 1 & 0 & 1 \\ 1 & 0 & 4 \\ 1 & 0 & 6 \end{pmatrix}.$$

The first column of  $X$  contains only 1's, while the second column of  $X$  contains a variable denoting the treatment group assignment and the third column contains the times of the repeated measurements.

# Square and Symmetric Matrices

A matrix is said to be *square* if it has the same number of rows and columns. A square matrix is *symmetric* if it equals its transpose. For example,

$$S = \begin{pmatrix} 2 & 3 & 7 & 11 \\ 3 & 9 & 1 & 2 \\ 7 & 1 & 5 & 8 \\ 11 & 2 & 8 & 4 \end{pmatrix}$$

is a symmetric matrix since it equals its transpose

$$S' = \begin{pmatrix} 2 & 3 & 7 & 11 \\ 3 & 9 & 1 & 2 \\ 7 & 1 & 5 & 8 \\ 11 & 2 & 8 & 4 \end{pmatrix}.$$

Examples of symmetric matrices that play an important role in the analysis of longitudinal data are the covariance and correlation matrices for the repeated measures on the same individuals.

Finally, a *diagonal* matrix is a special case of a symmetric square matrix that has non-zero elements only in the main diagonal positions, and zeros elsewhere. The main diagonal elements are those in the same row and column, from the upper left to the lower right corners of the matrix. For example,

$$D = \begin{pmatrix} 2 & 0 & 0 & 0 \\ 0 & 9 & 0 & 0 \\ 0 & 0 & 5 & 0 \\ 0 & 0 & 0 & 4 \end{pmatrix}$$

is a diagonal matrix. The diagonal matrix having all ones along the main diagonal is known as the *identity* matrix and is often denoted by  $I$  or  $I_n$ , where the subscript  $n$  denotes the dimension of the identity matrix. Thus

$$I_4 = \begin{pmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \end{pmatrix}.$$

# Arithmetic Operations

Addition and subtraction of matrices are defined only for matrices of the same dimension. That is, the matrices must share the same number of rows and the same number of columns. The sum of two matrices is obtained by adding their corresponding elements. For example,

$$\begin{pmatrix} 2 & 7 & 11 \\ 4 & 9 & 3 \\ 13 & 8 & 2 \end{pmatrix} + \begin{pmatrix} 3 & 2 & 14 \\ 7 & 8 & 4 \\ 6 & 5 & 9 \end{pmatrix} = \begin{pmatrix} 2+3 & 7+2 & 11+14 \\ 4+7 & 9+8 & 3+4 \\ 13+6 & 8+5 & 2+9 \end{pmatrix}$$
$$= \begin{pmatrix} 5 & 9 & 25 \\ 11 & 17 & 7 \\ 19 & 13 & 11 \end{pmatrix}.$$

Subtraction of matrices is defined in a similar way. For example,

$$\begin{pmatrix} 2 & 7 & 11 \\ 4 & 9 & 3 \\ 13 & 8 & 2 \end{pmatrix} - \begin{pmatrix} 3 & 2 & 14 \\ 7 & 8 & 4 \\ 6 & 5 & 9 \end{pmatrix} = \begin{pmatrix} 2-3 & 7-2 & 11-14 \\ 4-7 & 9-8 & 3-4 \\ 13-6 & 8-5 & 2-9 \end{pmatrix}$$
$$= \begin{pmatrix} -1 & 5 & -3 \\ -3 & 1 & -1 \\ 7 & 3 & -7 \end{pmatrix}.$$

# Scalar Multiplication of a Matrix

A scalar is a single number, as opposed to a vector or matrix of numbers. The scalar multiple of a matrix is formed by multiplying each element of the matrix by the scalar. For example, if

$$A = \begin{pmatrix} 2 & 7 & 11 & 5 \\ 4 & 9 & 3 & 1 \\ 13 & 8 & 2 & 6 \end{pmatrix}, \quad \text{then} \quad 2A = \begin{pmatrix} 4 & 14 & 22 & 10 \\ 8 & 18 & 6 & 2 \\ 26 & 16 & 4 & 12 \end{pmatrix}.$$

# Multiplication of Matrices

The multiplication of two matrices is somewhat more involved. The multiplication of two matrices  $A$  and  $B$ , denoted  $AB$ , is defined only if the number of columns of  $A$  is equal to the number of rows of  $B$ . For example, if  $A$  is a  $p \times q$  matrix and  $B$  is a  $q \times r$  matrix, then the product of the two matrices  $AB$  is a  $p \times r$  matrix. Letting  $C$  be the product of  $A$  and  $B$ ,

$$C = AB,$$

the  $(i, j)^{th}$  element of  $C$  is the sum of the products of the corresponding elements in the  $i^{th}$  row of  $A$  and the  $j^{th}$  column of  $B$ . Specifically, if  $c_{ij}$  is the element in the  $i^{th}$  row and the  $j^{th}$  column of the matrix  $C = AB$ , then

$$c_{ij} = \sum_{k=1}^q a_{ik}b_{kj} = a_{i1}b_{1j} + a_{i2}b_{2j} + \cdots + a_{iq}b_{qj}, \quad i = 1, \dots, p; \quad j = 1, \dots, r;$$

where  $q$  is the number of columns in  $A$  and the number of rows in  $B$ . Matrix multiplication is best understood by considering a simple example. Suppose

$$A = \begin{pmatrix} 2 & 7 & 11 \\ 4 & 9 & 3 \\ 13 & 8 & 2 \end{pmatrix} \quad \text{and} \quad B = \begin{pmatrix} 1 & 2 \\ 3 & 1 \\ 2 & 4 \end{pmatrix}$$

then

$$\begin{aligned} AB &= \begin{pmatrix} (2 \times 1) + (7 \times 3) + (11 \times 2) & (2 \times 2) + (7 \times 1) + (11 \times 4) \\ (4 \times 1) + (9 \times 3) + (3 \times 2) & (4 \times 2) + (9 \times 1) + (3 \times 4) \\ (13 \times 1) + (8 \times 3) + (2 \times 2) & (13 \times 2) + (8 \times 1) + (2 \times 4) \end{pmatrix} \\ &= \begin{pmatrix} 45 & 55 \\ 37 & 29 \\ 41 & 42 \end{pmatrix}. \end{aligned}$$

Note that the order of multiplication is very important. For example, if  $A$  and  $B$  are both square matrices of the same dimension, then  $AB$  is usually not equal to  $BA$ .

The multiplication of a vector by a matrix is a particularly important operation that plays a key role in longitudinal analysis. Let  $B$  be a  $p \times 1$  vector and  $X$  be a  $n \times p$  matrix. Then the product,

$$C = XB,$$

is a  $n \times 1$  vector with

$$c_i = \sum_{k=1}^p x_{ik}b_k, \quad i = 1, \dots, n;$$

where  $x_{ij}$  is the element in the  $i^{th}$  row and the  $j^{th}$  column of the matrix  $X$  and  $b_j$  is the element in the  $j^{th}$  row of the vector  $B$ . That is,

$$c_1 = x_{11}b_1 + x_{12}b_2 + \cdots + x_{1p}b_p,$$

$$c_2 = x_{21}b_1 + x_{22}b_2 + \cdots + x_{2p}b_p,$$

$$c_3 = x_{31}b_1 + x_{32}b_2 + \cdots + x_{3p}b_p,$$

and so on.

Let us return to the example introduced earlier, with repeated measures of blood lead levels obtained on four occasions. Letting  $Y$  denote the vector of repeated measurements of the response variable,

$$Y = \begin{pmatrix} Y_1 \\ Y_2 \\ Y_3 \\ Y_4 \end{pmatrix},$$

and  $X$  a matrix of covariates associated with the vector of repeated measurements,

$$X = \begin{pmatrix} X_{11} & X_{12} & X_{13} \\ X_{21} & X_{22} & X_{23} \\ X_{31} & X_{32} & X_{33} \\ X_{41} & X_{42} & X_{43} \end{pmatrix},$$

a linear regression model for the mean of each response can be expressed in vector and matrix notation as

$$E(Y) = X \beta,$$

where  $E(Y)$  denotes the expected value or mean of  $Y$  (see *Properties of Expectations and Variances* in Appendix B) and  $\beta$  is a  $3 \times 1$  vector of regression coefficients,

$$\beta = \begin{pmatrix} \beta_1 \\ \beta_2 \\ \beta_3 \end{pmatrix}.$$

Specifically, the product

$$E(Y) = X \beta,$$

is a  $4 \times 1$  vector

$$\begin{pmatrix} E(Y_1) \\ E(Y_2) \\ E(Y_3) \\ E(Y_4) \end{pmatrix} = \begin{pmatrix} X_{11} & X_{12} & X_{13} \\ X_{21} & X_{22} & X_{23} \\ X_{31} & X_{32} & X_{33} \\ X_{41} & X_{42} & X_{43} \end{pmatrix} \begin{pmatrix} \beta_1 \\ \beta_2 \\ \beta_3 \end{pmatrix} = \begin{pmatrix} \beta_1 X_{11} + \beta_2 X_{12} + \beta_3 X_{13} \\ \beta_1 X_{21} + \beta_2 X_{22} + \beta_3 X_{23} \\ \beta_1 X_{31} + \beta_2 X_{32} + \beta_3 X_{33} \\ \beta_1 X_{41} + \beta_2 X_{42} + \beta_3 X_{43} \end{pmatrix}.$$

That is,

$$E(Y) = X \beta,$$

is simply a shorthand representation for the following series of linear regression equations

$$E(Y_1) = \beta_1 X_{11} + \beta_2 X_{12} + \beta_3 X_{13},$$

$$E(Y_2) = \beta_1 X_{21} + \beta_2 X_{22} + \beta_3 X_{23},$$

$$E(Y_3) = \beta_1 X_{31} + \beta_2 X_{32} + \beta_3 X_{33},$$

$$E(Y_4) = \beta_1 X_{41} + \beta_2 X_{42} + \beta_3 X_{43}.$$

# Inverse

The *inverse* of a square matrix  $A$ , denoted  $A^{-1}$ , is defined as a square matrix whose elements are such that

$$AA^{-1} = A^{-1}A = I,$$

where  $I$  is the identity matrix, a diagonal matrix having all ones along the main diagonal. That is, the product of  $A$  by its inverse is equal to the identity matrix. The inverse of a square matrix does not always exist. The inverse of a matrix only exists if the matrix is *non-singular*.

In matrix algebra the inverse plays the role of the reciprocal, and thus multiplication by an inverse,  $A^{-1}$ , can loosely be thought of as “division” by the matrix  $A$ . Methods for calculating the inverse of a matrix will not be discussed here. In practice, the inverse of a matrix is usually obtained with the aid of a computer.

Finally, the *determinant* of a square matrix is a unique scalar (or single number) function of its elements and is denoted by  $|A|$ . For example, if  $A$  is a  $2 \times 2$  matrix,

$$A = \begin{pmatrix} a_{11} & a_{12} \\ a_{21} & a_{22} \end{pmatrix},$$

then the determinant of  $A$  is the following function of its elements

$$|A| = a_{11}a_{22} - a_{12}a_{21}.$$

The corresponding expression for the determinant of a  $3 \times 3$  square matrix, and of matrices of higher dimensions, is more involved and the details are not important. A useful property of the determinant is that it provides a test of whether the inverse of a matrix exists. In particular, if  $|A| \neq 0$ , then the inverse of  $A$  exists; if  $|A| = 0$ , then the matrix is said to be *singular* and the inverse of  $A$  does not exist.

The determinant also plays a role in the definition of the multivariate normal distribution (see Section 3.2). The multivariate normal density includes a term involving the determinant of the covariance matrix. The determinant of the covariance matrix is often referred to as the *generalized variance* and characterizes the salient features of the variation expressed by the covariance matrix in a single-number summary.

<sup>1</sup> In matrix algebra, vectors have a geometric meaning, denoting the coordinates of a point in Euclidean space. The geometric concept of the “length” (or magnitude) of a vector in Euclidean space has a very precise definition and technical meaning that is quite different from our informal use of the term here.

## ***Appendix B***

### ***Properties of Expectations and Variances***

Let  $Y$  denote a random variable that takes on values according to some probability density function if  $Y$  is continuous or some probability mass function if  $Y$  is discrete. The *expected value*, or *expectation*, of  $Y$  is simply its *mean* or average value and is usually denoted by

$$\mu = E(Y).$$

It is often referred to as the first *moment* of  $Y$ , since it describes the location of the center of the distribution. The precise definition of the expectation of  $Y$  is that it is a *weighted* average of all the possible values of  $Y$ , with weights determined by the probabilities associated with each possible value.

The *variance* of  $Y$ , often denoted by  $\sigma^2 = \text{Var}(F)$ , is a measure of the dispersion or variability around the mean or expected value of  $Y$ . The variance is often referred to as the second *central moment* of  $Y$  and is defined as

$$\sigma^2 = \text{Var}(Y) = E\{Y - E(Y)\}^2.$$

The variance is a weighted average of the squared deviations of  $Y$  around its mean. Because the variance is expressed in squared units of  $Y$ , a measure of variability in the original units of  $Y$  is given by the *standard deviation*

$$\sigma = \sqrt{\text{Var}(Y)}.$$

Finally, the covariance between two random variables,  $X$  and  $Y$ , is defined as

$$\text{Cov}(X, Y) = E[\{X - E(X)\}\{Y - E(Y)\}],$$

and is a measure of the *linear dependence* between  $X$  and  $Y$ . If  $X$  and  $Y$  are *independent*, then  $\text{Cov}(X, Y) = 0$ . Note that the covariance of a variable with itself is the variance,  $\text{Cov}(Y, Y) = \text{Var}(Y)$ .

# Properties of Expectations and Variances

Next we consider some properties of expectations and variances. Let  $X$  and  $Y$  be two (possibly dependent) random variables and let  $a$  and  $b$  denote non-random constants. Then the expectation operator,  $E(\cdot)$ , has the following five important properties:

1.  $E(a) = a$
2.  $E(bX) = b E(X)$
3.  $E(a + bX) = a + bE(X)$
4.  $E(aX + bY) = a E(X) + b E(Y)$
5.  $E(XY) \neq E(X) E(Y)$  (unless  $X$  and  $Y$  are *independent*)

Thus expectation is a linear operator in the sense that it respects or preserves the arithmetic operations of addition and multiplication by a constant. As a result the expected value of a linear function of  $Y$  (e.g.,  $a + bY$ ) is simply the same linear function of the expected value of  $Y$  (e.g.,  $a + bE(Y)$ ).

The variance operator,  $\text{Var}(\cdot)$ , has the following five important properties:

1.  $\text{Var}(a) = 0$
2.  $\text{Var}(bY) = b^2 \text{Var}(Y)$
3.  $\text{Var}(a + bY) = b^2 \text{Var}(Y)$
4.  $\text{Var}(aX + bY) = a^2 \text{Var}(X) + b^2 \text{Var}(Y) + 2ab \text{Cov}(X, Y)$
5.  $\text{Var}(aX - bY) = a^2 \text{Var}(X) + b^2 \text{Var}(Y) - 2ab \text{Cov}(X, Y)$

In particular, if  $X$  and  $Y$  are dependent, then

$$\text{Var}(X + Y) = \text{Var}(X) + \text{Var}(Y) + 2\text{Cov}(X, Y)$$

and

$$\text{Var}(X - Y) = \text{Var}(X) + \text{Var}(Y) - 2\text{Cov}(X, Y).$$

Finally, we note that the expectation and variance operators can also be applied to vectors of random variables. For example, let  $Y$  be a  $n \times 1$  (column) response vector (e.g., repeated measurements at  $n$  different occasions),

$$Y = \begin{pmatrix} Y_1 \\ Y_2 \\ \vdots \\ Y_n \end{pmatrix},$$

then

$$E(Y) = \begin{pmatrix} E(Y_1) \\ E(Y_2) \\ \vdots \\ E(Y_n) \end{pmatrix},$$

and

$$\text{Cov}(Y) = \begin{pmatrix} \text{Var}(Y_1) & \text{Cov}(Y_1, Y_2) & \dots & \text{Cov}(Y_1, Y_n) \\ \text{Cov}(Y_2, Y_1) & \text{Var}(Y_2) & \dots & \text{Cov}(Y_2, Y_n) \\ \vdots & \vdots & \ddots & \vdots \\ \text{Cov}(Y_n, Y_1) & \text{Cov}(Y_n, Y_2) & \dots & \text{Var}(Y_n) \end{pmatrix}.$$

## Appendix C

### *Critical Points for a 50:50 Mixture of Chi-Squared Distributions*

**Table C.1** Critical points for a 50:50 mixture of chi-squared distributions with  $q$  and  $q + 1$  degrees of freedom; right-hand tail probabilities.

$q$	Significance Level									
	0.3	0.2	0.1	0.05	0.025	0.01	0.005	0.0025	0.001	0.0005
0	0.28	0.71	1.64	2.71	3.84	5.41	6.63	7.88	9.55	10.83
1	1.76	2.50	3.81	5.14	6.48	8.27	9.63	11.00	12.81	14.18
2	3.06	3.98	5.53	7.05	8.54	10.50	11.97	13.43	15.36	16.80
3	4.29	5.36	7.09	8.76	10.38	12.48	14.04	15.59	17.61	19.13
4	5.49	6.68	8.57	10.37	12.10	14.32	15.97	17.59	19.69	21.27
5	6.66	7.96	10.00	11.91	13.74	16.07	17.79	19.47	21.66	23.29
6	7.82	9.21	11.38	13.40	15.32	17.76	19.54	21.29	23.55	25.23
7	8.97	10.44	12.74	14.85	16.86	19.38	21.23	23.04	25.37	27.10
8	10.10	11.66	14.07	16.27	18.35	20.97	22.88	24.74	27.13	28.91
9	11.23	12.87	15.38	17.67	19.82	22.52	24.49	26.40	28.86	30.68
10	12.35	14.06	16.67	19.04	21.27	24.05	26.07	28.02	30.54	32.40

*Source:* Adapted from Monette et al. (2002).

*Note:* Critical value  $c$  such that right-hand tail probability equals  $0.5 \times \Pr(\chi_q^2 > c) + 0.5 \times \Pr(\chi_{q+1}^2 > c)$ , where  $\chi_q^2$  and  $\chi_{q+1}^2$  denote chi-squared distributions with  $q$  and  $q + 1$  degrees of freedom, respectively.

## ***References***

- Adams, M.M., Wilson, H.G., Casto, D.L., Berg, C.J., McDermott, J.M., Gaudino, J.A., and McCarthy, B.J. (1997). Constructing reproductive histories by linking vital records. *American Journal of Epidemiology*, **145**, 339–348.
- Agresti, A. (2002). *Categorical Data Analysis*, 2nd ed. New York: Wiley.
- Agresti, A. (2010). *Analysis of Ordinal Categorical Data*, 2nd ed. New York: Wiley.
- Allison, P.D. (2005). *Fixed Effects Regression Methods for Longitudinal Data Using SAS*. Cary, NC: SAS Institute.
- Altham, P.M.E. (1984). Improving the precision of estimation by fitting a model. *Journal of the Royal Statistical Society, Series B*, **46**, 118–119.
- Altman, D.G. (1990). *Practical Statistics for Medical Research*. New York: Chapman and Hall/CRC Press.
- Amemiya, T. (1985). *Advanced Econometrics*. Cambridge, MA: Harvard University Press.
- Anderson, D.A. and Aitkin, M. (1985). Variance components models with binary response: Interviewer variability. *Journal of the Royal Statistical Society, Series B*, **47**, 203–210.
- Anderson, E.B. (1970). Asymptotic properties of conditional maximum likelihood estimators. *Journal of the Royal Statistical Society, Series B*, **32**, 283–301.
- Anderson, R.L. and Bancroft, T.A. (1952). *Statistical Theory in Research*. New York: McGraw-Hill.
- Bahadur, R.R. (1961). A representation of the joint distribution of responses to n dichotomous items. In Solomon, H. (editor), *Studies in Item Analysis and Prediction*, pp. 158–168. Palo Alto: Stanford University Press.
- Bandini, L.G., Must, A., Spadano, J.L. and Dietz, W.H. (2002). Relation of body composition, parental overweight, pubertal stage, and race-ethnicity to energy expenditure among premenarcheal girls. *American Journal of Clinical Nutrition*, **76**, 1040–1047.
- Basagana, X. and Spiegelman, D. (2010). Power and sample size calculations for longitudinal studies comparing rates of change with a time-varying exposure. *Statistics in Medicine*, **29**, 181–192.
- Becker, M.P. and Balagtas, C.C. (1993). Marginal modeling of binary cross-over data. *Biometrics*, **49**, 997–1009.
- Begg, M.D. and Parides, M.K. (2003). Separation of individual-level and cluster-level covariate effects in regression analysis of correlated data. *Statistics in Medicine*, **22**, 2591–2602.
- Berlin, J.A., Kimmel, S.E., Ten Have, T.R. and Sammel, M.D. (1999). An empirical comparison of several clustered data approaches under confounding due to cluster effects in the analysis of complications of coronary angioplasty. *Biometrics*, **55**, 470–476.
- Bock, R.D. (1975). *Multivariate Statistical Methods in Behavioral Research*. New York: McGraw-Hill.
- Bombardier, C., Ware, J., Russell, I.J., Larson, M., Chalmers, A. and Read, J.L. (1986). Auranofin therapy and quality of life in patients with rheumatoid arthritis. Results of a multicenter trial. *American Journal of Medicine*, **81**, 565–578.
- Box, G.E.P. (1950). Problems in the analysis of growth and wear data. *Biometrics*, **6**, 362–389.
- Breslow, N.E. (2005). Whither PQL? In Lin, D.Y. and Heagerty, P.J. (editors), *Proceedings of the Second Seattle Symposium in Biostatistics: Analysis of Correlated Data*. New York: Springer.
- Breslow, N.E. and Clayton, D.G. (1993). Approximate inference in generalized linear mixed models. *Journal of the American Statistical Association*, **88**, 9–25.

- Breslow, N.E. and Lin, X. (1995). Bias correction in generalized linear models with a single component of dispersion. *Biometrika*, **82**, 81–91.
- Brown, H. and Prescott, R. (1999). *Applied Mixed Models in Medicine*. Chichester, UK: Wiley.
- Browne, W.J., Draper, D., Goldstein, H. and Rasbash, J. (2002). Bayesian and likelihood methods for fitting multilevel models with complex level-1 variation. *Computational Statistics and Data Analysis*, **39**, 203–225.
- Brumback, B. and Rice, J. (1998). Smoothing spline models for the analysis of nested and crossed samples of curves. *Journal of the American Statistical Association*, **93**, 961–993.
- Burton, P., Gurrin, L. and Sly, P. (1998). Tutorial in biostatistics: Extending the simple linear regression model to account for correlated responses: An introduction to generalized estimating equations and multi-level modelling. *Statistics in Medicine*, **17**, 1261–1291.
- Byar, D. and Blackard, C. (1977). Comparisons of placebo, pyridoxine, and topical thiotepa in preventing recurrence of stage I bladder cancer. *Urology*, **10**, 556–561.
- Cameron, A.C. and Trivedi, P.K. (1998). *Regression Analysis of Count Data*. Cambridge, UK: Cambridge University Press.
- Carlin, B.P. and Louis, T.A. (2000). *Bayes and Empirical Bayes Methods for Data Analysis*, 2nd ed. New York: Chapman and Hall/CRC Press.
- Carroll, R.J., Hall, P., Apanasovich, T.V. and Lin, X. (2004). Histospline method in nonparametric regression models with application to clustered/longitudinal data. *Statistica Sinica*, **14**, 649–674.
- Chi, E.M. and Reinsel, G.C. (1989). Models of longitudinal data with random effects and AR(1) errors. *Journal of the American Statistical Association*, **84**, 452–459.
- Cnaan, A., Laird, N.M. and Slasor, P. (1997). Tutorial in biostatistics: Using the general linear mixed model to analyze unbalanced repeated measures and longitudinal data. *Statistics in Medicine*, **16**, 2349–2380.
- Cochran, W.G. and Cox, G.M. (1957). *Experimental Designs*, 2nd ed. New York: Wiley.
- Cole, J.W.L. and Grizzle, J.E. (1966). Applications of multivariate analysis of variance to repeated measurements experiments. *Biometrics*, **22**, 810–828.
- Collett, D. (1991). *Modelling Binary Data*. London: Chapman and Hall/CRC Press.
- Conaway, M.R. (1992). The analysis of repeated categorical measurements subject to nonignorable nonresponse. *Journal of the American Statistical Association*, **87**, 817–824.
- Cook, R.D. (1977). Detection of influential observations in linear regression. *Technometrics*, **19**, 15–18.
- Cook, R.D. and Weisberg, S. (1982). *Residuals and Influence in Regression*. New York: Chapman and Hall/CRC Press.
- Cook, R.J., Zeng, L. and Yi, G.Y. (2004). Marginal analysis of incomplete longitudinal binary data: A cautionary note on LOCF imputation. *Biometrics*, **60**, 820–828.
- Cox, D.R. (1961). Tests of separate families of hypotheses. In *Proceedings of the Fourth Berkeley Symposium on Mathematical Statistics and Probability*, Vol. 1, pp. 105–123. Berkeley, CA: University of California Press.
- Cox, D.R. (1970). *Analysis of Binary Data*, 1st ed. New York: Chapman and Hall/CRC Press.
- Cox, D.R. (1972). Regression models and life tables (with discussion). *Journal of the Royal Statistical Society, Series B*, **34**, 187–220.
- Cox, D.R. and Wermuth, N. (1996). *Multivariate Dependencies. Models, Analysis and Interpretation*. New York: Chapman and Hall/CRC Press.
- Crainiceanu, C. and Ruppert, D. (2004). Restricted likelihood ratio tests for longitudinal models. *Statistica Sinica*, **14**, 713–729.
- Crowder, M.J. and Hand, D.J. (1990). *Analysis of Repeated Measures*. London: Chapman and

Hall/CRC Press.

- Dale, J.R. (1984). Local versus global association for bivariate ordered responses. *Biometrika*, **71**, 507–514.
- Danford, M.B., Hughes, H.M. and McNee, R.C. (1960). On the analysis of repeated measurements experiments. *Biometrics*, **16**, 547–565.
- Daniels, M. and Gatsonis, C. (1999). Hierarchical generalized linear models in the analysis of variations in health care utilization. *Journal of the American Statistical Association*, **94**, 29–42.
- Daniels, M.J. and Hogan, J.W. (2008). *Missing Data in Longitudinal Studies: Strategies for Bayesian Modeling and Sensitivity Analysis*. New York: Chapman and Hall/CRC Press.
- Davidian, M. and Giltinan, D.M. (1995). *Nonlinear Models for Repeated Measurement Data*. London: Chapman and Hall/CRC Press.
- Davis, C.S. (1991). Semi-parametric and non-parametric methods for the analysis of repeated measurements with applications to clinical trials. *Statistics in Medicine*, **10**, 1959–1980.
- Davis, C.S. (2002) *Statistical Methods for the Analysis of Repeated Measurements*. New York: Springer.
- Dawber, T.R. (1980). *The Framingham Study: The Epidemiology of Atherosclerotic Disease*. Cambridge, MA: Harvard University Press.
- Dawber, T.R., Meadors, G.F. and Moore, F.E.J. (1951). Epidemiological approaches to heart disease: The Framingham Study. *American Journal of Public Health*, **41**, 279–286.
- De Backer, M., De Vroey, C., Lesaffre, E., Scheyns, I., and De Keyser, P. (1998). Twelve weeks of continuous oral therapy for toenail onychomycosis caused by dermatophytes: A double-blind comparative trial of terbinafine 250 mg/day versus itraconazole 200 mg/day. *Journal of the American Academy of Dermatology*, **38**, 57–63.
- Dempster, A.P., Laird, N.M. and Rubin, D.B. (1977). Maximum likelihood from incomplete data via the EM algorithm. *Journal of the Royal Statistical Society, Series B*, **39**, 1–38.
- Diggle, P.J. (1988). An approach to the analysis of repeated measurements. *Biometrics*, **44**, 959–971.
- Diggle, P.J. (1990). *Time Series: A Biostatistical Introduction*. Oxford, UK: Oxford University Press.
- Diggle, P.J., Heagerty, P., Liang, K.-Y. and Zeger, S.L. (2002). *Analysis of Longitudinal Data*, 2nd ed. New York: Oxford University Press.
- Dobson, A.J. (1990). *An Introduction to Generalized Linear Models*. London: Chapman and Hall/CRC Press.
- Dockery, D.W., Berkey, C.S., Ware, J.H., Speizer, F.E. and Ferris, B.G. (1983). Distribution of FVC and FEV<sub>1</sub> in children 6 to 11 years old. *American Review of Respiratory Disease*, **128**, 405–412.
- Donnelly, C.A., Laird, N.M. and Ware, J.H. (1995). Prediction and creation of smooth curves for temporally correlated longitudinal data. *Journal of the American Statistical Association*, **90**, 984–989.
- Drum, M. and McCullagh, P. (1993). Comment on “Regression models for discrete longitudinal responses”. *Statistical Science*, **8**, 300–301.
- Dufouil, C., Brayne, C. and Clayton, D. (2004). Analysis of longitudinal studies with death and drop-out: A case study. *Statistics in Medicine*, **23**, 2215–2226.
- Durbán, M., Harezlak, J., Wand, M.P. and Carroll, R.J. (2005). Simple fitting of subject-specific curves for longitudinal data. *Statistics in Medicine*, **24**, 1153–1167.
- Durbin, J. (1954). Errors in variables. *Review of the International Statistical Institute*, **22**, 23–32.
- Durbin, J. (1960). Estimation of parameters in time-series regression models. *Biometrika*, **47**, 139–153.

- Efron, B. and Morris, C. (1977). Stein's paradox in statistics. *Scientific American*, **236**, 119–127.
- Ekholm, A. (1991). Fitting regression models to a multivariate binary response. In Rosenqvist, G., Juselius, K., Nordstrom, K. and Palmgren, J. (editors), *A Spectrum of Statistical Thought: Essays in Statistical Theory, Economics, and Population Genetics in Honour of Johan Fellman*, pp. 19–32. Helsingfors: Swedish School of Economics and Business Administration.
- Ekholm, A., Smith, P.W.F. and McDonald, J.W. (1995). Marginal regression analysis of a multivariate binary response. *Biometrika*, **82**, 847–854.
- Emrich, L.J. and Piedmonte, M.R. (1992). On some small sample properties of generalized estimating equation estimates for multivariate dichotomous outcomes. *Journal of Statistical Computation and Simulation*, **41**, 19–29.
- Everitt, B.S. (1995). The analysis of repeated measures: A practical review with examples. *Statistician*, **44**, 113–135.
- Fahrmeir, L. and Tutz, G. (2001). *Multivariate Statistical Modelling Based on-Generalized Linear Models*, 2nd ed. New York: Springer.
- Fay, M.P. and Graubard, B.I. (2001). Small-sample adjustments for Wald-type tests using sandwich estimators. *Biometrics*, **57**, 1198–1206.
- Fay, M.P., Graubard, B.I., Freedman, L.S. and Midthune, D.N. (1998). Conditional logistic regression with sandwich estimators: Application to a metaanalysis. *Biometrics*, **54**, 195–208.
- Feldman, H.A. (1988). Families of lines: Random effects in linear regression analysis. *Journal of Applied Physiology*, **64**, 1721–1732.
- Firth, D. (1991). Generalized linear models. In Hinkley, D.V., Reid, N. and Snell, E.J. (editors), *Statistical Theory and Modelling. In Honour of Sir David Cox, FRS*. London: Chapman and Hall/CRC Press.
- Fisher, R.A. (1925). *Statistical Methods for Research Workers*. Edinburgh, UK: Oliver and Boyd.
- Fitzmaurice, G.M. (2001). A conundrum in the analysis of change. *Nutrition*, **17**, 360–361.
- Fitzmaurice, G., Davidian, M., Verbeke, G. and Molenberghs, G. (2009). *Longitudinal Data Analysis*. Boca Raton, FL: Chapman and Hall/CRC Press.
- Fitzmaurice, G.M. and Laird, N.M. (1993). A likelihood-based method for analysing longitudinal binary responses. *Biometrika*, **80**, 141–151.
- Fitzmaurice, G.M., Laird, N.M. and Rotnitzky, A.G. (1993). Regression models for discrete longitudinal responses (with discussion). *Statistical Science*, **8**, 248–309.
- Fitzmaurice, G.M., Laird, N.M., Zahner, G.E.P. and Daskalakis, C. (1995). Bivariate logistic regression analysis of childhood psychopathology ratings using multiple informants. *American Journal of Epidemiology*, **142**, 1194–1203.
- Flay, B.R., Miller, T.Q., Hedeker, D., Siddiqui, O., Brannon, B.R., Johnson, C.A., Hansen, W.B., Sussman, S. and Dent, C. (1995). The television, school and family smoking prevention and cessation project: VIII. Student outcomes and mediating variables. *Preventive Medicine*, **24**, 29–410.
- Fracom, S.F., Chuang-Stein, C. and Landis, J.R. (1989). A log-linear model for ordinal data to characterize differential change among treatments. *Statistics in Medicine*, **8**, 571–582.
- Freund, R.J., Littell, R.C. and Spector, P.C. (1986). *SAS System for Linear Models*. Cary, NC: SAS Institute.
- Friedman, G.D., Cutter, G.R., Donahue, R., Hughes, G.H., Hulley, S., Jacobs, D.R., Liu, K. and Savage, P.J. (1988). CARDIA: Study design, recruitment, and some characteristics of the examined subjects. *Journal of Clinical Epidemiology*, **41**, 1105–1116.
- Gardiner, J.C., Luo, Z. and Roman, L.A. (2009). Fixed effects, random effects and GEE: What are the differences? *Statistics in Medicine*, **28**, 221–239.
- Geisser, S. (1963). Multivariate analysis of variance for a special covariance case. *Journal of the*

- American Statistical Association*, **58**, 660–669.
- Gelman, A., Carlin, J.B., Stern, H.S. and Rubin, D.B. (2003). *Bayesian Data Analysis*, 2nd ed. New York: Chapman and Hall/CRC Press.
- Gelman, A. and Hill, J. (2007). *Data Analysis using Regression and Multilevel/Hierarchical Models*. Cambridge, UK: Cambridge University Press.
- Ghiden, W., Lesaffre, E. and Eilers, P. (2004). Smooth random effects distribution in a linear mixed model. *Biometrics*, **60**, 945–953.
- Gibbons, R.D. and Hedeker, D. (1994). Application of random-effects probit regression models. *Journal of Consulting and Clinical Psychology*, **62**, 285–296.
- Gibbons, R.D., Hedeker, D., Waternaux, C. and Davis, J.M. (1988). Random regression models: A comprehensive approach to the analysis of longitudinal psychiatric data. *Psychopharmacological Bulletin*, **24**, 438–443.
- Gill, J. (2000). *Generalized Linear Models: A Unified Approach*. Thousand Oaks, CA: Sage Publications.
- Gilmour, A.R., Anderson, R.D. and Rae, A.L. (1985). The analysis of binomial data by generalized linear mixed models. *Biometrika*, **72**, 593–599.
- Glonek, G.F.V. (1996). A class of regression models for multivariate categorical responses. *Biometrika*, **83**, 15–28.
- Glonek, G.F.V. and McCullagh, P. (1995). Multivariate logistic models. *Journal of the Royal Statistical Society, Series B*, **57**, 533–546.
- Glymour, M.M., Weuve, J., Berkman, L.F., Kawachi, I. and Robins, J.M. (2005). When is baseline adjustment useful in analyses of change: An example with education and cognitive change. *American Journal of Epidemiology*, **162**, 267–278.
- Godambe, V.P. (1960). An optimum property of regular maximum likelihood estimation. *Annals of Mathematical Statistics*, **31**, 1208–1212.
- Goldstein, H. (2003). *Multilevel Statistical Methods*, 3rd ed. London: Edward Arnold.
- Goldstein, H., Browne, W. and Rasbach, J. (2002). Multilevel modelling of medical data. *Statistics in Medicine*, **21**, 3291–3315.
- Goldwasser, M. and Fitzmaurice, G.M. (2001). Multivariate linear regression analysis of childhood psychopathology using multiple informant data. *International Journal of Methods in Psychiatric Research*, **10**, 1–10.
- Gourieroux, C., Monfort, A. and Trognon, A. (1984). Pseudo-maximum likelihood methods: Theory. *Econometrica*, **52**, 681–700.
- Graubard, B.I. and Korn, E.L. (1994). Regression analysis with clustered data. *Statistics in Medicine*, **13**, 509–522.
- Green, P.J. (1985). Linear models for field trials, smoothing and cross-validation. *Biometrika*, **72**, 523–537.
- Greenberg, E.R., Baron, J.A., Stevens, M.M., Stukel, T.A., Mandel, J.S., Spencer, S.K., Elias, P.M., Lowe, N., Nierenberg, D.W., Bayrd, G. and Vance, J.C. (1989). The Skin Cancer Prevention Study: Design of a clinical trial of beta-carotene among persons at high risk for nonmelanoma skin cancer. *Controlled Clinical Trials*, **10**, 153–166.
- Greenberg, E.R., Baron, J.A., Stukel, T.A., Stevens, M.M., Mandel, J.S., Spencer, S.K., Elias, P.M., Lowe, N., Nierenberg, D.W., Bayrd, G., Vance, J.C., Freeman, D.H., Clendenning, W.E., Kwan, T. and the Skin Cancer Prevention Study Group (1990). A clinical trial of beta carotene to prevent basal-cell and squamous-cell cancers of the skin. *New England Journal of Medicine*, **323**, 789–795.
- Greenhouse, S.W. and Geisser, S. (1959). On methods in the analysis of profile data. *Psychometrika*, **32**, 95–112.
- Grizzle, J.E. and Allen, D.W. (1969). Analysis of growth and dose response curves. *Biometrics*, **25**,

- Grizzle, J.E., Starmer, C.F. and Koch, G.G. (1969). Analysis of categorical data by linear models. *Biometrics*, **15**, 489–504.
- Gunsolley, J.C., Getchell, C. and Chinchilli, V.M. (1995). Small sample characteristics of generalized estimating equations. *Communications in Statistics, Simulation and Computation*, **24**, 869–878.
- Gurrin, L.C., Scurrah, K.J. and Hazelton, M.L. (2005). Tutorial in biostatistics: Spline smoothing with linear mixed models. *Statistics in Medicine*, **24**, 3361–3381.
- Hand, D.J. and Crowder, M.J. (1996). *Practical Longitudinal Data Analysis*. New York: Chapman and Hall/CRC Press.
- Hand, D.J. and Taylor, C.C. (1987). *Multivariate Analysis of Variance and Repeated Measures*. New York: Chapman and Hall/CRC Press.
- Harezlak, J., Naumova, E. and Laird, N.M. (2007). LongCrisp: A test for bump hunting in longitudinal data. *Statistics in Medicine*, **26**, 1383–1397.
- Harezlak, J., Ryan, L.M., Giedd, J.N. and Lange, N. (2005). Individual and population penalized regression splines for accelerated longitudinal designs. *Biometrics*, **61**, 1037–1048.
- Hartley, H.O. and Rao, J.N.K. (1967). Maximum-likelihood estimation for the mixed analysis of variance model. *Biometrika*, **54**, 93–108.
- Harville, D.A. (1974). Bayesian inference for variance components using only error contrasts. *Biometrika*, **61**, 383–385.
- Harville, D.A. (1977). Maximum likelihood approaches to variance component estimation and to related problems. *Journal of the American Statistical Association*, **72**, 320–338.
- Hastie, T.J. and Tibshirani, R.J. (1990). *Generalized Additive Models*. London: Chapman and Hall.
- Hausman, J.A. (1978). Specification tests in econometrics. *Econometrica*, **46**, 1251–1271.
- Hedeker, D. and Gibbons, R.D. (1994). A random-effects ordinal regression model for multilevel analysis. *Biometrics*, **50**, 933–944.
- Hedeker, D. and Gibbons, R.D. (1996). MIXOR: A computer program for mixed-effects ordinal regression analysis. *Computer Methods and Programs in Biomedicine*, **49**, 157–176.
- Hedeker, D. and Gibbons, R.D. (2006). *Longitudinal Data Analysis*. New York: Wiley.
- Hedeker, D., Gibbons, R.D. and Waternaux, C. (1999). Sample size estimation for longitudinal designs with attrition: Comparing time-related contrasts between two groups. *Journal of Educational and Behavioral Statistics*, **24**, 70–93.
- Henderson, C.R. (1963). Selection index and expected genetic advance. In Hanson, W.D. and Robinson, H.F. (editors), *Statistical Genetics and Plant Breeding*. Washington, DC: National Academy of Sciences-National Research Council.
- Henry, K., Erice, A., Tierney, C., Balfour, H.H. Jr, Fischl, M.A., Kmack, A., Liou, S.H., Kenton, A., Hirsch, M.S., Phair, J., Martinez, A. and Kahn J.O. for the AIDS Clinical Trial Group 193A Study Team (1998). A randomized, controlled, double-blind study comparing the survival benefit of four different reverse transcriptase inhibitor therapies (three-drug, two-drug, and alternating drug) for the treatment of advanced AIDS. *Journal of Acquired Immune Deficiency Syndromes and Human Retrovirology*, **19**, 339–349.
- Hernández, B., Gortmaker, S.L., Laird, N.M., Colditz, G.A., Parra-Cabrera, S. and Peterson, K.E. (2000). Validity and reproducibility of a physical activity and inactivity questionnaire for Mexico City's schoolchildren. *Salud Publico de Mexico*, **42**, 315–323.
- Heyting, A., Tolboom, J. and Essers, J. (1992). Statistical handling of dropouts in longitudinal clinical trials. *Statistics in Medicine*, **11**, 2043–2061.
- Hilbe, J.M. (2007). *Negative Binomial Regression*. Cambridge, UK: Cambridge University Press.
- Hinkley, D.V. and Wang, S. (1991). Efficiency of robust standard errors for regression coefficients.

*Communications in Statistics, Theory and Methods*, **20**, 1–11.

Hogan, J.W. and Laird, N.M. (1996). Intention to treat analysis for incomplete repeated measures data. *Biometrics*, **52**, 1002–1017.

Hogan, J.W., Roy, J. and Korkontzelou, C. (2004). Tutorial in biostatistics: Handling drop-out in longitudinal studies. *Statistics in Medicine*, **23**, 1455–1497.

Horton, N.J. and Lipsitz, S.R. (2001). Multiple imputation in practice: Comparison of software packages for regression models with missing variables. *American Statistician*, **55**, 244–254.

Horvitz, D.G. and Thompson, D.J. (1952). A generalization of sampling without replacement from a finite universe. *Journal of the American Statistical Association*, **47**, 663–685.

Hosmer, D.W. Jr. and Lemeshow, S. (2000). *Applied Logistic Regression*, 2nd ed. New York: Wiley.

Hox, J. (2002). *Multilevel Analysis: Techniques and Applications*. London: Lawrence Erlbaum Associates.

Hubbard, A.E., Ahern, J., Fleischer, N.L., Van der Laan, M., Lippman, S.A., Jewell, N., Bruckner, T. and Satariano, W.A. (2010). To GEE or not to GEE: Comparing population average and mixed models for estimating the associations between neighborhood risk factors and health. *Epidemiology*, **21**, 467–474.

Huber, P.J. (1967). The behavior of maximum likelihood estimates under non-standard conditions. In *Proceedings of the Fifth Berkeley Symposium on Mathematical Statistics and Probability*. Vol 1, pp. 221–233. Berkeley, CA: University of California Press.

Huber, P.J. (1981). *Robust Statistics*. New York: Wiley.

Javaras, K.N. and Fitzmaurice, G.M. (2009). Semiparametric regression modeling of incomplete longitudinal outcomes: Stratifying on informative missingness predictors. *Statistics in Biopharmaceutical Research*, **1**, 48–65.

Jennrich, R.I. and Schluchter, M.D. (1986). Unbalanced repeated-measures models with structured covariance matrices. *Biometrics*, **42**, 805–820.

Johnson, R.A. and Wichern, D.W. (2002). *Applied Multivariate Statistical Analysis*, 5th ed. Englewood Cliff, NJ: Prentice-Hall.

Jones, B. and Donev, A.N. (1996). Modelling and design of cross-over trials. *Statistics in Medicine*, **15**, 1435–1446.

Jones, B. and Kenward, M.G. (1989). *Design and Analysis of Cross-over Trials*. London: Chapman and Hall/CRC Press.

Jung, S.-H. and Ahn, C.W. (2005). Sample size for a two-group comparison of repeated binary measurements using GEE. *Statistics in Medicine*, **24**, 2583–2596.

Kackar, R.N. and Harville, D.A. (1981). Unbiasedness of two-stage estimation and prediction procedures for mixed linear models. *Communications in Statistics, Theory and Methods*, **10**, 1249–1261.

Kakwani, N.C. (1967). The unbiasedness of Zellner's seemingly unrelated regression equations estimators. *Journal of the American Statistical Association*, **62**, 141–142.

Kauermann, G. and Carroll, R.J. (2001). The sandwich variance estimator: Efficiency properties and coverage probability of confidence intervals. *Journal of the American Statistical Association*, **96**, 1387–1396.

Kenward, M.G. and Molenberghs, G. (1999). Parametric models for incomplete continuous and categorical longitudinal data. *Statistical Methods in Medical Research*, **8**, 51–83.

Kenward, M.G. and Molenberghs, G. (2009). Last observation carried forward: A crystal ball? *Journal of Biopharmaceutical Statistics*, **19**, 872–888.

Kenward, M.G. and Roger, J.H. (1997). Small sample inference for fixed effects from restricted maximum likelihood. *Biometrics*, **53**, 983–997.

Keselman, H.J. and Keselman, J.C. (1984). The analysis of repeated measures designs in medical

- research. *Statistics in Medicine*, **3**, 185–195.
- Khatri, C.G. (1966). A note on a MANOVA model applied to problems in growth curves. *Annals of the Institute of Statistical Mathematics*, **18**, 75–86.
- Kleinbaum, D.G., Kupper, L.L., Muller, K.E. and Nizam, A. (1998). *Applied Regression Analysis and Multivariable Methods*, 3rd ed. Pacific Grove, CA: Duxbury Press.
- Koch, G.G., Amara, I.A., Brown, B.W., Colton, T. and Gillings, D.B. (1989). A two-period crossover design for the comparison of two active treatments and placebo. *Statistics in Medicine*, **8**, 487–504.
- Koch, G.G., Carr, G.J., Amara, I.A., Stokes, M.E. and Uryniak, T.J. (1990). Categorical data analysis. In Berry, D.A. (editor), *Statistical Methodology in the Pharmaceutical Sciences*. New York: Marcel Dekker.
- Koch, G.G., Landis, J.R., Freeman, J.L., Freeman, D.H. and Lehnen, R.G. (1977). A general methodology for the analysis of experiments with repeated measurement of categorical data. *Biometrics*, **33**, 133–158.
- Kramer, M.S., Chalmers, B., Hodnett, E.D., Sevkovskaya, Z., Dzikovich, I. and Shapiro, S. for the PROBIT Study Group (2001). Promotion of breastfeeding intervention trial (PROBIT): A randomized trial in the Republic of Belarus. *Journal of the American Medical Association*, **285**, 413–420.
- Kreft, I.I. and de Leeuw, J. (1998). *Introducing Multilevel Modeling*. London: Sage Publications.
- Laird, N.M. (1983). Further comparative analyses of pretest-posttest research designs. *American Statistician*, **37**, 329–330.
- Laird, N.M. (1988). Missing data in longitudinal studies. *Statistics in Medicine*, **7**, 305–315.
- Laird, N.M., Donnelly, C. and Ware, J.H. (1992). Longitudinal models with continuous responses. *Statistical Methods in Medical Research*, **1**, 225–247.
- Laird, N.M., Lange, N. and Stram, D.O. (1987). Maximum likelihood computations with repeated measures: Application of the EM algorithm. *Journal of the American Statistical Association*, **82**, 97–105.
- Laird, N.M., Skinner, J. and Kenward, M.G. (1992). An analysis of two-period crossover designs with carry-over effects. *Statistics in Medicine*, **11**, 1967–1979.
- Laird, N.M. and Ware, J.H. (1982). Random effects models for longitudinal data. *Biometrics*, **38**, 963–974.
- Lang, J.B. and Agresti, A. (1994). Simultaneously modeling joint and marginal distributions of multivariate categorical responses. *Journal of the American Statistical Association*, **89**, 625–632.
- Langford, I.H., Bentham, G. and McDonald, A. (1998). Multilevel modelling of geographically aggregated health data: A case study on malignant melanoma mortality and UV exposure in the European community. *Statistics in Medicine*, **17**, 41–58.
- Lauer, R.M., Clarke, W.R. and Burns, T.L. (1997). Obesity in childhood: The Muscatine Study. *Acta Padiatrica Scandinavica*, **38**, 432–437.
- Leppik, I., Dreifuss, F.E., Porter, R., Bowman, T., Santilli, N., Jacobs, M., Crosby, C., Cloyd, J., Stackman, J., Sutula, T.P., Graves, N., Welty, T., Vickery, T., Bundage, R., Gates, J., Gummit, R. and Gutierrez, A. (1987). A controlled study of progabide in partial seizure: Methodology and results. *Neurology*, **37**, 963–968.
- Lesaffre, E. and Spiessens, B. (2001). On the effect of the number of quadrature points in a logistic random-effects model: An example. *Applied Statistics*, **50**, 325–335.
- Leyland, A.H. and Goldstein, H. (2001). *Multilevel Modelling of Health Statistics*. Chichester, UK: Wiley.
- Liang, K.-Y. and Zeger, S.L. (1986). Longitudinal data analysis using generalized linear models. *Biometrika*, **73**, 13–22.

- Liang, K.-Y. and Zeger, S.L. (1995). Inference based on estimating functions in the presence of nuisance parameters (with discussion). *Statistical Science*, **10**, 158–199.
- Liang, K.-Y., Zeger, S.L. and Qaqish, B. (1992). Multivariate regression analyses for categorical data (with discussion). *Journal of the Royal Statistical Society, Series B*, **54**, 2–24.
- Lin, D.Y. and Wei, L.J. (1989). The robust inference for the Cox proportional hazards model. *Journal of the American Statistical Association*, **84**, 1074–1078.
- Lin, D.Y., Wei, L.J. and Ying, Z. (2002). Model-checking techniques based on cumulative residuals. *Biometrics*, **58**, 1–12.
- Lin, X. and Breslow, N.E. (1996). Bias correction in generalized linear mixed models with multiple components of dispersion. *Journal of the American Statistical Association*, **91**, 1007–1017.
- Lindley, D.V. and Smith, A.F.M. (1972). Bayes estimates for the linear model. *Journal of the Royal Statistical Society, Series B*, **34**, 1–41.
- Lindsey, J.K. (1999). *Models for Repeated Measurements*, 2nd ed. Oxford, UK: Clarendon Press.
- Lindstrom, M.J. and Bates, D.M. (1988). Newton-Raphson and EM algorithms for linear mixed effects models for repeated measures data. *Journal of the American Statistical Association*, **83**, 1014–1022.
- Lipsitz, S.R. and Fitzmaurice, G.M. (1994). Sample size for repeated measures studies with binary responses. *Statistics in Medicine*, **13**, 1233–1239.
- Lipsitz, S.R., Laird, N.M. and Harrington, D.P. (1990). Maximum likelihood regression methods for paired binary data. *Statistics in Medicine*, **9**, 1417–1425.
- Littell, R.C., Milliken, G.A., Stroup, W.W. and Wolfinger, R.D. (1996). *SAS System for Mixed Models*. Cary, NC: SAS Institute.
- Littell, R.C., Pendergast, J. and Natarajan, R. (2000). Tutorial in biostatistics: Modelling covariance structure in the analysis of repeated measures data. *Statistics in Medicine*, **19**, 1793–1819.
- Little, R.J.A. (1995). Modelling the drop-out mechanism in repeated-measures studies. *Journal of the American Statistical Association*, **90**, 1112–1121.
- Little, R.J.A. and Rubin, D.B. (2001). *Statistical Analysis with Missing Data*, 2nd ed. New York: Wiley.
- Little, R.J. and Yau, L. (1996). Intent-to-treat analysis for longitudinal studies with drop-outs. *Biometrics*, **52**, 1324–1333.
- Liu, G. and Liang, K.-Y. (1997). Sample size calculations for studies with correlated observations. *Biometrics*, **53**, 937–947.
- Liu, I. and Agresti, A. (2005). The analysis of ordered categorical data: An overview and a survey of recent developments. *TEST*, **14**, 1–73.
- Longford, N. (1993). *Random Coefficient Models*. Oxford, UK: Oxford University Press.
- Lord, F.M. (1967). A paradox in the interpretation of group comparisons. *Psychological Bulletin*, **68**, 304–305.
- Lorr, M. and Klett, C.J. (1966). *Inpatient Multidimensional Psychiatric Scale: Manual*, (rev.). Palo Alto, CA: Consulting Psychologists Press.
- Louis, T.A., Robins, J., Dockery, D.W., Spiro, A. and Ware, J.H. (1986). Explaining discrepancies between longitudinal and cross-sectional models. *Journal of Chronic Diseases*, **39**, 831–839.
- Machin, D., Farley, T., Busca, B., Campbell, M. and d'Arcangues, C. (1988). Assessing changes in vaginal bleeding patterns in contracepting women. *Contraception*, **38**, 165–179.
- Maddala, G.S. (1971). The use of variance components models in pooling cross section and time series data. *Econometrica*, **39**, 341–358.
- Mancl, L.A. and DeRouen, T.A. (2001). A covariance estimator for GEE with improved small-sample properties. *Biometrics*, **57**, 126–134.

- Matthews, J.N.S. (1994). Multi-period crossover trials. *Statistical Methods in Medical Research*, **3**, 383–405.
- Matthews, J.N.S., Altman, D.G., Campbell, M.J. and Royston, P. (1990). Analysis of serial measurements in medical research. *British Medical Journal*, **300**, 230–235.
- McCullagh, P. and Nelder, J. (1989). *Generalized Linear Models*, 2nd ed. London: Chapman and Hall/CRC Press.
- McCulloch, C.E. (2005). Repeated measures ANOVA, R.I.P.? *Chance*, **18**, 28–33.
- McCulloch, C.E. and Searle, S.R. (2001). *Generalized, Linear, and Mixed Models*. New York: Wiley.
- Miller, J.J. (1977). Asymptotic properties of maximum likelihood estimates in the mixed model of the analysis of variance. *Annals of Statistics*, **5**, 746–762.
- Molenberghs, G. and Kenward, M.G. (2007). *Missing Data in Clinical Studies*. Chichester: Wiley.
- Molenberghs, G. and Lesaffre, E. (1994). Marginal modelling of correlated ordinal data using a multivariate Plackett distribution. *Journal of the American Statistical Association*, **89**, 633–644.
- Molenberghs, G. and Ritter, L. (1996). Likelihood and quasi-likelihood based methods for analyzing multivariate categorical data, with the association between outcome of interest. *Biometrics*, **52**, 1121–1133.
- Molenberghs, G., Thijs, H., Jansen, I., Beunckens, C., Kenward, M.G., Mallinckrodt, C. and Carroll, R.J. (2004). Analyzing incomplete longitudinal clinical trial data. *Biostatistics*, **5**, 445–464.
- Molenberghs, G. and Verbeke, G. (2005). *Models for Discrete Longitudinal Data*. New York: Springer.
- Molenberghs, G., Verbeke, G., Demétrio, C.G.B. and Vieira, A.M.C. (2010). A family of generalized linear models for repeated measures with normal and conjugate random effects. *Statistical Science*, **25**, 325–347.
- Monette, G., Kwan, E., Rivilis, A. and Shao, Q. (2002). A first look at multilevel models. Unpublished manuscript. Department of Mathematics and Statistics, York University.
- Morrison, D.F. (1972). The analysis of a single sample of repeated measurements. *Biometrics*, **28**, 55–71.
- Morrison, D.F. (1990). *Multivariate Statistical Methods*, 3rd ed. New York: McGraw-Hill.
- Muller, K.E., LaVange, L.M., Ramey, S.L., and Ramey, C.T. (1992). Power calculations for general linear multivariate models including repeated measures applications. *Journal of the American Statistical Association*, **87**, 1209–1226.
- Myers, R.H., Montgomery, D.C. and Vining, G.G. (2001). *Generalized Linear Models: With Applications in Engineering and the Sciences*. New York: Wiley.
- National Center for Health Statistics (1992). Sample design: Third National Health and Nutrition Examination Survey. *Vital and Health Statistics*, Series 2, No. 113.
- National Center for Health Statistics (1994). Plan and operation of the third National Health and Nutrition Examination Survey, 1988–94. *Vital and Health Statistics*, Series 1, No. 32.
- Naumova, E.N., Must, A. and Laird, N.M. (2001). Evaluating the impact of “critical periods” in longitudinal studies of growth using piecewise mixed effects models. *International Journal of Epidemiology*, **30**, 1332–1341.
- Nelder, J.A. and Wedderburn, R.W.M. (1972). Generalized linear models. *Journal of the Royal Statistical Society, Series A*, **135**, 370–384.
- Neter, J., Kutner, M.H., Nachtsheim, C.J. and Wasserman, W. (1996). *Applied Linear Regression Models*, 3rd ed. Homewood, IL: Richard D. Irwin.
- Neuhaus, J.M. (2001). Assessing change with longitudinal and clustered binary data. *Annual Review of Public Health*, **22**, 115–128.

- Neuhaus, J.M. and Kalbfleisch, J.D. (1998). Between- and within-cluster covariate effects in the analysis of clustered data. *Biometrics*, **54**, 638–645.
- Neuhaus, J.M., Kalbfleisch, J.D. and Hauck, W.W. (1991). A comparison of cluster-specific and population-averaged approaches for analyzing correlated binary data. *International Statistical Review*, **59**, 25–35.
- Newey, W. and McFadden, D. (1994). Large sample estimation and hypothesis testing. In Engle, R.F. and McFadden, D. (editors), *Handbook of Econometrics*, Vol 4. Amsterdam: North-Holland.
- Neyman, J. and Scott, E.L. (1948). Consistent estimates based on partially consistent observations. *Econometrica*, **16**, 1–32.
- O'Connell, M.A. and Wolfinger, R.D. (1997). Spatial regression models, response surfaces, and process optimization. *Journal of Computational and Graphical Statistics*, **6**, 224–241.
- Omar, R.Z., Wright, E.M., Turner, R.M. and Thompson, S.G. (1999). Analyzing repeated measures data: A practical comparison of methods. *Statistics in Medicine*, **18**, 1587–1603.
- Overall, J.E. and Doyle, S.R. (1994). Estimating sample sizes for repeated measurement designs. *Controlled Clinical Trials*, **15**, 100–123.
- Pagano, M. and Gauvreau, K. (2000). *Principles of Biostatistics*, 2nd ed. Pacific Grove, CA: Duxbury Press.
- Pan, W. (2001). Sample size and power calculations with correlated binary data. *Controlled Clinical Trials*, **22**, 211–227.
- Pan, W. (2002). A note on the use of marginal likelihood and conditional likelihood in analyzing clustered data. *American Statistician*, **56**, 171–174.
- Pan, W. and Connett, J.E. (2002). Selecting the working correlation structure in generalized estimating equations with application to the lung health study. *Statistica Sinica*, **12**, 475–490.
- Pan, W., Louis, T.A. and Connett, J.E. (2000). A note on marginal linear regression with correlated response data. *American Statistician*, **54**, 191–195.
- Patterson, H.D. and Thompson, R. (1971). Recovery of inter-block information when block sizes are unequal. *Biometrika*, **58**, 545–554.
- Pepe, M.S. and Anderson, G.A. (1994). A cautionary note on inference for marginal regression models with longitudinal data and general correlated response data. *Communications in Statistics, Simulation and Computation*, **23**, 939–951.
- Phillips, S.M., Bandini, L.G., Compton, D.V., Naumova, E.N. and Must, A. (2003). A longitudinal comparison of body composition by total body water and bioelectrical impedance in adolescent girls. *Journal of Nutrition*, **133**, 1419–1425.
- Pierce, D.A. and Sands, B.R. (1975). Extra-Bernoulli variation in binary data. Technical Report 46, Department of Statistics, Oregon State University.
- Pinheiro, J.C. and Bates, D.M. (1995). Approximations to the log-likelihood function in the nonlinear mixed-effects model. *Journal of Computational and Graphical Statistics*, **4**, 12–35.
- Potthoff, R.F. and Roy, S.N. (1964). A generalized multivariate analysis of variance model useful especially for growth curve problems. *Biometrika*, **51**, 313–326.
- Pourahmadi, M. (1999). Joint mean-covariance models with applications to longitudinal data: Unconstrained parameterization. *Biometrika*, **86**, 677–690.
- Preisser, J.S., Lohman, K.K. and Rathouz, P.J. (2002). Performance of weighted estimating equations for longitudinal binary data with dropouts missing at random. *Statistics in Medicine*, **21**, 3035–3054.
- Preisser, J.S. and Qaqish, B.F. (1996). Deletion diagnostics for generalised estimating equations. *Biometrika*, **83**, 551–562.
- Price, C.J., Kimmel, C.A., Tyl, R.W. and Marr, M.C. (1985). The developmental toxicity of ethylene glycol in rats and mice. *Toxicological Applications in Pharmacology*, **81**, 113–127.

- Qaqish, B.F. and Liang, K.-Y. (1992). Marginal models for correlated binary responses with multiple classes and multiple levels of nesting. *Biometrics*, **48**, 939–950.
- Rabe-Hesketh, S. and Skrondal, A. (2008). *Multilevel and Longitudinal Modeling Using Stata*, 2nd ed. College Station, TX: Stata Press.
- Raghunathan, T.E., Lepkowski, J.M., van Hoewyk, J. and Solenberger, P. (2001). A multivariate technique for multiply imputing missing values using a sequence of regression models. *Survey Methodology*, **27**, 85–95.
- Rao, C.R. (1958). Some statistical methods for comparison of growth curves. *Biometrics*, **14**, 1–17.
- Rao, C.R. (1965). The theory of least squares when the parameters are stochastic and its application to the analysis of growth curves. *Biometrika*, **52**, 447–458.
- Rathouz, P.J. (2004). Fixed effects models for longitudinal binary data with drop-outs missing at random. *Statistica Sinica*, **14**, 969–988.
- Raudenbush, S.W. and Bryk, A.S. (2002). *Hierarchical Linear Models: Applications and Data Analysis Methods*, 2nd ed. Newbury Park, CA: Sage Publications.
- Raudenbush, S.W. and Liu, X.F. (2001). Effects of study duration, frequency of observation, and sample size on power in studies of group differences in polynomial change. *Psychological Methods*, **6**, 387–401.
- Reilly, M. and Pepe, M. (1997). The relationship between hot-deck multiple imputation and weighted likelihood. *Statistics in Medicine*, **16**, 5–19.
- Rencher, A.C. (2002). *Methods of Multivariate Analysis*, 2nd ed. New York: Wiley.
- Rijcken, B., Schouten, J.P., Weiss, S.T., Speizer, F.E. and van der Lende, R. (1987). The relationship of nonspecific bronchial responsiveness to respiratory symptoms in a random population sample. *American Review of Respiratory Disease*, **136**, 62–68.
- Robins, J.M., Greenland, S. and Hu, F.-C. (1999). Estimation of the causal effect of a time-varying exposure on the marginal mean of a repeated binary outcome (with discussion). *Journal of the American Statistical Association*, **94**, 687–712.
- Robins, J.M., Rotnitzky, A. and Zhao, L.P. (1995). Analysis of semiparametric regression models for repeated outcomes in the presence of missing data. *Journal of the American Statistical Association*, **90**, 106–121.
- Robinson, G.K. (1991). That BLUP is a good thing: The estimation of random effects (with discussion). *Statistical Science*, **6**, 15–51.
- Rochon J. (1998). Application of GEE procedures for sample size calculations in repeated measures experiments. *Statistics in Medicine*, **17**, 1643–1658.
- Rogan, W.J., Dietrich, K.N., Ware, J.H., Dockery, D.W., Salganik, M., Radcliffe, J., Jones, R.L., Ragan, N.B., Chisolm, J.J. and Rhoads, G.G. (2001). The effect of chelation therapy with succimer on neuropsychological development in children exposed to lead. *New England Journal of Medicine*, **344**, 1421–1426.
- Rosenbaum, P.R. and Rubin, D.B. (1983). Assessing sensitivity to an unobserved binary covariate in an observational study with binary outcome. *Journal of the Royal Statistical Society, Series B*, **45**, 212–218.
- Rosenman, R.H., Brand, R.J., Jenkins, C.D., Friedman, M., Straus, R. and Wurm, M. (1975). Coronary heart disease in the Western Collaborative Study: Final follow-up experience of 8  $\frac{1}{2}$  years. *Journal of the American Medical Association*, **233**, 872–877.
- Rowell, J.G. and Walters, D.E. (1976). Analysing data with repeated observations on each experimental unit. *Journal of Agricultural Science*, **87**, 423–432.
- Royall, R.M. (1986). Model robust confidence intervals using maximum likelihood estimators. *International Statistical Review*, **54**, 221–226.
- Rubin, D.B. (1976). Inference and missing data. *Biometrika*, **63**, 581–592.

- Rubin, D.B. (1978). Multiple imputations in sample surveys: A phenomeno-logical Bayesian approach to nonresponse. In *Proceedings of the Survey Research Methods Section*, pp. 20–34. Washington, DC: American Statistical Association.
- Rubin, D.B. (1987). *Multiple Imputation for Nonresponse in Surveys*. New York: Wiley.
- Rubin, D.B. (1996). Multiple imputation after 18+ years. *Journal of the American Statistical Association*, **91**, 473–489.
- Rubin, D.B. and Schenker, N. (1991). Multiple imputation in health-care databases: an overview and some applications. *Statistics in Medicine*, **10**, 585–598.
- Ruppert, D., Wand, M.P. and Carroll, R.J. (2003). *Semiparametric Regression*. Cambridge, UK: Cambridge University Press.
- Ruppert, D., Wand, M.P. and Carroll, R.J. (2009). Semiparametric regression during 2003–2007. *Electronic Journal of Statistics*, **3**, 1193–1256.
- Russell, T.S. and Bradley, R.A. (1958). One-way variances in a two-way classification. *Biometrika*, **45**, 111–129.
- Satterthwaite, F.E. (1946). An approximate distribution of estimates of variance components. *Biometrics Bulletin*, **2**, 110–114.
- Schafer, J.L. (1997). *Analysis of Incomplete Multivariate Data*. London: Chapman and Hall/CRC Press.
- Schafer, J.L. (1999). Multiple imputation: A primer. *Statistical Methods in Medical Research*, **8**, 3–15.
- Schall, R. (1991). Estimation in generalized linear models with random effects. *Biometrika*, **78**, 719–727.
- Scheffé, H. (1959). *The Analysis of Variance*. New York: Wiley.
- Schildcrout, J.S. and Heagerty, P.J. (2005). Regression analysis of longitudinal binary data with time-dependent environmental covariates: Bias and efficiency. *Biostatistics*, **6**, 633–652.
- Schlesselman, J.J. (1973a). Planning a longitudinal study: I. Sample size determination. *Journal of Chronic Diseases*, **26**, 553–560.
- Schlesselman, J.J. (1973b). Planning a longitudinal study: II. Frequency of measurement and study duration. *Journal of Chronic Diseases*, **26**, 561–570.
- Schoenfield, L.J., Lachin, J.M., Baum, R.A., Habig, R.L., Hanson, R.F., Hersh, T., Hightower, N.C., Hofmann, A.F., Lasser, E.C., Marks, J.W., Mekhjian, H., Okun, R., Schaefer, R.A., Shaw, L., Soloway, R.D., Thistle, J.L., Thomas, F.B. and Tyor, M.P. (1981). National Cooperative Gallstone Study: A controlled trial of the efficacy and safety of chenodeoxycholic acid for dissolution of gallstones. *Annals of Internal Medicine*, **95**, 257–282.
- Schwartz, D. and Lellouch, J. (1967). Explanatory and pragmatic attitudes in therapeutical trials. *Journal of Chronic Diseases*, **20**, 637–648.
- Scott, A.J. and Holt, D. (1982). The effect of two-stage sampling on ordinary least squares methods. *Journal of the American Statistical Association*, **77**, 848–854.
- Self, S.G. and Liang, K.-Y. (1987). Asymptotic properties of maximum likelihood estimators and likelihood ratio tests under nonstandard conditions. *Journal of the American Statistical Association*, **82**, 605–610.
- Selvin, S. (1995). *Practical Biostatistical Methods*. Belmont, CA: Duxbury Press.
- Senn, S. (2002). *Cross-over Trials in Clinical Research*, 2nd ed. New York: Wiley.
- Silvapulle, M.J. (1996). A test in the presence of nuisance parameters. *Journal of the American Statistical Association*, **91**, 1690–1693.
- Silvapulle, M.J. and Silvapulle, P. (1995). A score test against one-sided alternatives. *Journal of the American Statistical Association*, **90**, 342–349.

- Singer, J.D. (1998). Using SAS PROC MIXED to fit multilevel models, hierarchical models, and individual growth models. *Journal of Educational and Behavioral Statistics*, **25**, 323–355.
- Singer, J.D. and Willett, J.B. (2003). *Applied Longitudinal Data Analysis: Modeling Change and Event Occurrence*. New York: Oxford University Press.
- Skellam, J.G. (1948). A probability distribution derived from the binomial distribution by regarding the probability of success as variable between the sets of trials. *Journal of the Royal Statistical Society, Series B*, **10**, 257–261.
- Snedecor, G.W. and Cochran, W.G. (1967). *Statistical Methods*, 6th ed. Ames, IA: Iowa State Press.
- Snijders, T.A.B. and Bosker, R.J. (1993). Standard errors and sample sizes for two-level research. *Journal of Educational Statistics*, **18**, 237–259.
- Snijders, T.A.B. and Bosker, R.J. (1999). *Multilevel Analysis: An Introduction to Basic and Advanced Multilevel Modeling*. London: Sage Publications.
- Sochett, E.B., Daneman, D., Clarson, C. and Ehrlich, R.M. (1987). Factors affecting and patterns of residual insulin secretion during the first year of type 1 (insulin-dependent) diabetes mellitus in children. *Diabetologia*, **30**, 453–459.
- Speed, T. (1991). Comment on “That BLUP is a good thing: The estimation of random effects,” by G.K. Robinson. *Statistical Science*, **6**, 42–44.
- Spiegelhalter, D.J., Thomas, A. and Best, N.G. (1999). *WinBUGS Version 1.2 User Manual*. MRC Biostatistics Unit, Cambridge, UK.
- Stein, C. (1955). Inadmissibility of the usual estimator for the mean of a multivariate normal distribution. In *Proceedings of the Third Berkeley Symposium on Mathematical Statistics and Probability*, Vol. 1, pp. 197–206. Berkeley, CA: University of California Press.
- Stiratelli, R., Laird, N.M. and Ware, J.H. (1984). Random effects models for serial observations with binary response. *Biometrics*, **40**, 961–971.
- Stokes, M.E., Davis, C.S. and Koch, G.G. (1995). *Categorical Data Analysis using the SAS System*. Cary, NC: SAS Institute.
- Stram, D.O. and Lee, J.W. (1994). Variance components testing in the longitudinal mixed effects model. *Biometrics*, **50**, 1171–1177.
- Stram, D.O. and Lee, J.W. (1995). Correction to “Variance components testing in the longitudinal mixed effects model.” *Biometrics*, **51**, 1196.
- Stukel, T.A. (1993). Comparison of methods for the analysis of longitudinal interval count data. *Statistics in Medicine*, **12**, 1339–1351.
- Subramanian, S.V., Jones, K. and Duncan, C. (2003). Multilevel methods for public health research. In Kawachi, I. and Berkman, L.F. (editors), *Neighborhoods and Health*. New York: Oxford University Press.
- Sullivan, L.M., Dukes, K.A. and Losina, E. (1999). Tutorial in biostatistics: An introduction to hierarchical linear modeling. *Statistics in Medicine*, **18**, 855–888.
- Tchetgen, E.J. and Coull, B.A. (2006). A diagnostic test for the mixing distribution in a generalized linear mixed model. *Biometrika*, **93**, 1003–1010.
- Thall, P.F. and Vail, S.C. (1990). Some covariance models for longitudinal count data with overdispersion. *Biometrics*, **46**, 657–671.
- Thompson, R. (1985). Comment on “Some aspects of the spline smoothing approach to non-parametric regression curve fitting,” by B.W. Silverman. *Journal of the Royal Statistical Society, Series B*, **47**, 43–44.
- Thompson, W.A. (1962). The problem of negative estimates of variance components. *Annals of Mathematical Statistics*, **33**, 273–289.
- Treatment of Lead-Exposed Children (TLC) Trial Group (2000). Safety and efficacy of succimer in toddlers with blood lead levels of 20–44 µg/dL. *Pediatric Research*, **48**, 593–599.

- Tsiatis, A.A. (2006). *Semiparametric Theory and Missing Data*. New York: Springer.
- van Buuren, S. (2007). Multiple imputation of discrete and continuous data by fully conditional specification. *Statistical Methods in Medical Research*, **16**, 219–242.
- van Buuren, S., Brand, J.P.L., Groothuis-Oudshoorn, C.G.M. and Rubin, D.B. (2006). Fully conditional specification in multivariate imputation. *Journal of Statistical Computation and Simulation*, **76**, 1049–1064.
- van der Lende, R., Kok, T.J., Peset, R., Quanjer, P.H., Schouten, J.P. and Orie, N.G.M. (1981). Decreases in VC and FEV<sub>1</sub> with time: Indicators for effects of smoking and air pollution. *Bulletin of European Physiopathology and Respiration*, **17**, 775–792.
- Van Marter, L.J., Leviton, A., Kuban, K.C.K., Pagano, M. and Allred, E.N. (1990). Maternal glucocorticoid therapy and reduced risk of bronchopulmonary dysplasia. *Pediatrics*, **86**, 331–336.
- Verbeke, G. and Molenberghs, G. (2000). *Linear Mixed Models for Longitudinal Data*. New York: Springer.
- Verbeke, G. and Molenberghs, G. (2003). The use of score tests for inference on variance components. *Biometrics*, **59**, 254–262.
- Wahba, G. (1978). Improper priors, spline smoothing and the problem of guarding against model errors in regression. *Journal of the Royal Statistical Society, Series B*, **40**, 364–372.
- Wang, N., Carroll, R.J. and Lin, X. (2005). Efficient semiparametric marginal estimation for longitudinal/clustered data. *Journal of the American Statistical Association*, **100**, 147–157.
- Wang, Y. (1998). Mixed effects smoothing spline analysis of variance. *Journal of the Royal Statistical Society, Series B*, **60**, 159–174.
- Ware, J.H. (1985). Linear models for the analysis of longitudinal studies. *American Statistician*, **39**, 95–101.
- Ware, J.H. (2003). Interpreting incomplete data in studies of diet and weight loss. *New England Journal of Medicine*, **348**, 2136–2137.
- Ware, J.H., Dockery, D., Louis, T.A., Xu, X., Ferris, B.G. and Speizer, F.E. (1990). Longitudinal and cross-sectional estimates of pulmonary function decline in never-smoking adults. *American Journal of Epidemiology*, **132**, 685–700.
- Ware, J.H. and Liang, K.-Y. (1996). The design and analysis of longitudinal studies: A historical perspective. In Armitage, P. and David, H.A. (editors), *Advances in Biometry*. New York: Wiley.
- Waternaux, C., Laird, N.M. and Ware, J.H. (1989). Methods for analysis of longitudinal data: Blood-lead concentration and cognitive development. *Journal of the American Statistical Association*, **84**, 33–41.
- Waternaux, C. and Ware, J.H. (1991). Unconditional linear models for analysis of longitudinal data. In Dwyer, J.H., Feinleib, M., Lippert, P. and Hoffmeister, H. (editors), *Statistical Models for Longitudinal Studies of Health*. New York: Oxford University Press.
- Wei, L.J. and Lachin, J.M. (1984). Two-sample asymptotically distribution-free tests for incomplete multivariate observations. *Journal of the American Statistical Association*, **79**, 653–661.
- Wei, L.J., Lin, D.Y. and Weissfeld, L. (1989). Regression analysis of multivariate incomplete failure time data by modeling marginal distributions. *Journal of the American Statistical Association*, **84**, 1065–1073.
- Weisberg, S. (1985). *Applied Linear Regression*, 2nd ed. New York: Wiley.
- Welham, S.J., Cullis, B.R., Kenward, M.G. and Thompson, R. (2006). The analysis of longitudinal data using mixed model L-splines. *Biometrics*, **62**, 392–401.
- Welsh, A. H., Lin, X. and Carroll, R.J. (2002). Marginal longitudinal nonparametric regression: Locality and efficiency of spline and kernel methods. *Journal of the American Statistical Association*, **97**, 482–493.
- White, H. (1982). Maximum likelihood estimation of misspecified models. *Econometrica*, **50**, 1–26.

- White, I.R., Horton, N.J., Carpenter, J. and Pocock, S.J. (2011). Strategy for intention to treat analysis in randomised trials with missing outcome data. *British Medical Journal*, **342**, 910–912.
- Williams, D.A. (1982). Extra-binomial variation in logistic linear models. *Applied Statistics*, **31**, 144–148.
- Winer, B.J. (1971). *Statistical Principles in Experimental Design*, 2nd ed. New York: McGraw-Hill.
- Wishart, J. (1938). Growth-rate determinations in nutrition studies with the bacon pig, and their analysis. *Biometrika*, **30**, 16–28.
- Wolfinger, R. (1993). Laplace's approximation for nonlinear mixed models. *Biometrika*, **80**, 791–795.
- Wong, G.Y. and Mason, W.M. (1985). The hierarchical logistic regression model for multilevel analysis. *Journal of the American Statistical Association*, **80**, 513–524.
- Woolson, R.F. and Clarke, W.R. (1984). Analysis of categorical incomplete longitudinal data. *Journal of the Royal Statistical Society, Series A*, **147**, 87–99.
- Wu, D.-M. (1973). Alternative tests of independence between stochastic regressors and disturbances. *Econometrica*, **41**, 733–750.
- Yates, F. (1935). Complex experiments (with discussion). *Supplement to the Journal of the Royal Statistical Society*, **2**, 181–247.
- Zahner, G.E.P., Jacobs, J.H., Freeman, D.H. and Trainor, K.F. (1993). Rural-urban child psychopathology in a northeastern U.S. state: 1986–1989. *Journal of the American Academy of Child and Adolescent Psychiatry*, **32**, 378–387.
- Zahner, G.E.P., Pawelkiewicz, W., DeFrancesco, J.J. and Adnopo, J. (1992). Children's mental health service needs and utilization patterns in an urban community: An epidemiological assessment. *Journal of the American Academy of Child and Adolescent Psychiatry*, **31**, 951–960.
- Zeger, S.L. and Karim, M.R. (1991). Generalized linear models with random effects: A Gibbs sampling approach. *Journal of the American Statistical Association*, **86**, 79–86.
- Zeger, S.L. and Liang, K.-Y (1986). Longitudinal data analysis for discrete and continuous outcomes. *Biometrics*, **42**, 121–130.
- Zeger, S.L., Liang, K.-Y. and Albert, P.S. (1988). Models for longitudinal data: A generalized estimating equation approach. *Biometrics*, **44**, 1049–1060.
- Zhao, L.P., Prentice, R. and Self, S. (1992). Multivariate mean parameter estimation by using a partly exponential model. *Journal of the Royal Statistical Society, Series B*, **54**, 805–811.
- Zimmerman, D.L. and Nunez-Anton, V. (2001). Parametric modelling of growth curve data: An overview (with discussion). *Test*, **10**, 1–73.

# ***Index***

Adjustment for baseline response  
    alternative methods  
    analysis of covariance (ANCOVA)  
    Lord's paradox

AIDS Clinical Trial Group (ACTG) Study 193A

Akaike information criterion (AIC)

Alachua County Study of Mental Impairment

Analysis of covariance (ANCOVA)

Analysis of response profiles  
    design matrix  
    general linear model formulation  
    hypotheses  
    missing data  
    single degree of freedom tests  
    strengths and weaknesses

Analysis of variance (ANOVA)  
    repeated measures

Area under the curve (AUC)

Arthritis Clinical Trial

Autoregressive covariance

Balanced design

Balanced incomplete block design

Banded covariance

Baseline response  
    adjustment for  
    Lord's paradox

Bayesian information criterion (BIC)

Bernoulli distribution

Best linear unbiased predictor (BLUP)

Between-subject variability

Binomial distribution

BLUP

Broken-stick model

Canonical link function

CARDIA Study

Carryover effects

Centering

Change scores

Cholesky decomposition

Cholesky factorization

Clinical Trial of an Anti-epileptic Drug

Clinical Trial of Antibiotics for Leprosy

Clinical Trial of Contracepting Women

Clinical Trial of Patients with Insomnia

Clinical Trial of Patients with Respiratory Illness

Cluster  
Clustered data  
Clustering  
Cluster-randomized trial  
Cohort effects  
Compositional covariate  
Compound symmetry  
Conditional logistic regression  
Conditional maximum likelihood  
  conditional logistic regression  
  McNemar's test  
Conditional mean  
Connecticut Child Surveys  
Consistent estimator  
Contrast matrix  
Contrasting marginal and mixed effects models  
Correlation  
  constraints with binary data  
  definition  
  implications for longitudinal data  
  intra-cluster correlation  
Correlation matrix  
Covariance  
  Cholesky decomposition  
  compound symmetry  
  conditional  
  definition  
  interdependence on model for the mean  
  marginal  
  modeling of  
  nested models  
  non-standard likelihood theory  
  positive-definite  
  random effects structure  
  semi-variogram  
  unstructured  
Covariance matrix  
Covariance pattern models  
  autoregressive  
  banded  
  choice among  
  compound symmetry  
  exponential  
  hybrid model  
  strengths and weaknesses  
  Toeplitz  
Crossover design  
  carryover effects

Crossover Study of Pain Relief for Tension Headache

Crossover Trial on Cerebrovascular Deficiency

Cross-sectional effects

Cross-sectional study

defining feature

design

Cumulative sum of residuals

Decomposition of between- and within-subject effects

Denominator degrees of freedom

Kenward and Roger method

Satterthwaite approximation

Design issues

Design matrix

Developmental Toxicity Study of Ethylene Glycol

Difference scores

Dispersion parameter

Dropout

available-data analysis

baseline value carried forward

complete-case analysis

imputation

inverse probability weighted methods

last observation carried forward (LOCF)

last value carried forward (LVCF)

multiple imputation

pattern mixture model

predictive mean matching

propensity score methods

regression method

weighting methods

worst value carried forward

Eastern Connecticut Child Survey (ECCS)

*See also* Connecticut Child Surveys

Effect size

EM algorithm

Empirical Bayes estimator

Empirical variance estimator

*See also* “Sandwich” variance estimator

Endogenous covariate

Estimation

conditional maximum likelihood

generalized least squares (GLS)

marginal quasi-likelihood (MQL)

maximum likelihood (ML)

ordinary least squares (OLS)

penalized quasi-likelihood (PQL)

restricted maximum likelihood (REML)

## Examples

AIDS Clinical Trial Group (ACTG) Study 193A  
Alachua County Study of Mental Impairment  
Arthritis Clinical Trial  
Clinical Trial of an Anti-epileptic Drug  
Clinical Trial of Antibiotics for Leprosy  
Clinical Trial of Contracepting Women  
Clinical Trial of Patients with Insomnia  
Clinical Trial of Patients with Respiratory Illness  
Connecticut Child Surveys  
Crossover Study of Pain Relief for Tension Headache  
Crossover Trial on Cerebrovascular Deficiency  
Developmental Toxicity Study of Ethylene Glycol  
Exercise Therapy Trial  
Log C-Peptide Concentration in Children with Diabetes  
Malignant Melanoma Mortality and Ultraviolet Light Exposure  
MIT Growth and Development Study  
Muscatine Coronary Risk Factor (MCRF) Study  
National Cooperative Gallstone Study (NCGS)  
Onychomycosis Study  
Six Cities Study of Air Pollution and Health  
Skin Cancer Prevention Study  
Study of Bladder Cancer Tumors  
Study of Dental Growth  
Study of Low Birth Weight Infants  
Study of Progesterone Metabolite Concentration  
Study of Risk Factors for Coronary Heart Disease (CHD)  
Study of Weight Gain  
Television, School and Family Smoking Prevention and Cessation Project  
Treatment of Lead-Exposed Children (TLC) Trial  
U.S. Centers for Disease Control Study of Infant Birth Weight  
U.S. National Institute of Mental Health (NIMH) Schizophrenia Collaborative Study  
Vlagtwedde–Vlaardingen Study  
Exercise Therapy Trial  
Exogeneity  
Exogenous covariate  
Expectation  
  properties  
Exponential covariance  
Exponential family of distributions  
External covariate  
Extra-binomial variation  
  *See also* Overdispersion  
Extra-Poisson variation  
  *See also* Overdispersion  
First-order autoregression  
Fisher-scoring algorithm  
Fitted value

Fixed effects

Framingham Heart Study

Gamma distribution

Gaussian quadrature

General linear model

Generalized estimating equations (GEE)

algorithm

empirical variance estimator

properties of

sandwich variance estimator

working covariance matrix

Generalized least squares (GLS)

properties

Generalized linear mixed effects models

approximate methods of estimation

conditional maximum likelihood

contrasting marginal and mixed effects models

estimation and inference

interpretation of parameters

marginal quasi-likelihood

overdispersion

penalized quasi-likelihood

Generalized linear models

canonical link function

dispersion parameter

distributional assumption

estimation

extensions to longitudinal data

linear predictor

link function

ordinal data

overdispersion

overview

salient features

scale parameter

systematic component

variance function

Generalized variance

Group-randomized trial

Growth curve models

*See also* Linear mixed effects models

Hausman test

Hierarchical data

*See also* Multilevel data

Hierarchical models

*See also* Multilevel models

Imputation

Incidental parameters problem

Incomplete data

*See also* Missing data

Independence

Inference

likelihood ratio test

model parameters

multivariate Wald test

Wald test

Informative dropout

Informative missingness

Internal covariate

Intra-cluster correlation

Inverse probability weighted methods

Laplace approximation

Leverage

Likelihood function

Likelihood ratio test

Linear fixed effects models

choosing between fixed and random effects models

fixed versus random effects

Linear mixed effects model

decomposition of between- and within-subject effects

Linear mixed effects models

choosing between fixed and random effects models

conditional mean

fixed versus random effects

marginal mean

NIH method

population-averaged mean

prediction

random effects covariance structure

random intercept model

shrinkage

subject-specific mean

two-stage formulation

Linear regression

Litter effects

Locally-weighted regression

Logistic distribution

Logistic regression

example

overdispersion

underlying latent variable

Log-linear regression

example

offset

overdispersion

Longitudinal data

basic concepts

consequences of ignoring correlation

correlation

dependence

descriptive methods of analysis

distributional assumptions

historical approaches to analysis

notation

objectives of analysis

sources of correlation

Longitudinal effects

Longitudinal study

balanced design

defining feature

design issues

primary goal

unbalanced design

Lord's paradox

Lowess

Mahalanobis distance

Malignant Melanoma Mortality and Ultraviolet Light Exposure

MANOVA

Marginal likelihood

Marginal mean

Marginal models

contrasting marginal and mixed effects models

distributional assumptions

estimation

generalized estimating equations (GEE)

illustrative examples

implicit assumption

interpretation of model parameters

residual diagnostics

specification of

Marginal quasi-likelihood

basis of approximation

relation to generalized estimating equations (GEE)

Markov chain Monte Carlo (MCMC)

Matrices

arithmetic operations

basic concepts

definitions

determinant

inverse

square

symmetric

transpose

Maximal model

Maximum likelihood estimation

likelihood function

maximum likelihood estimate (MLE)

score function

McNemar's test

Mean response

analysis of response profiles

broken-stick model

linear splines

linear trends

maximal model

modeling of

parametric curves

piecewise linear trend

polynomial trends

quadratic trends

residual diagnostics

saturated model

semiparametric curves

truncated line function

Measurement error

Missing at random (MAR)

Missing completely at random (MCAR)

Missing data

available-data analysis

baseline value carried forward

complete-case analysis

dropout

EM algorithm

implications for analysis

imputation

inverse probability weighted methods

last observation carried forward (LOCF)

last value carried forward (LVCF)

monotone missing data pattern

multiple imputation

pattern mixture model

predictive mean matching

propensity score methods

regression method

weighting methods

worst value carried forward

Missing data mechanisms

covariate-dependent missingness

hierarchy of

ignorable

informative

- missing at random (MAR)
- missing completely at random (MCAR)
- nonignorable
- not missing at random (NMAR)
- response indicator variables

Mistimed measurements

MIT Growth and Development Study

Modeling the covariance

Moving average

Multilevel data

Multilevel generalized linear models

Multilevel linear models

Multilevel models

- estimation

- three-level generalized linear models

- three-level linear models

- two-level generalized linear models

- two-level linear models

Multiple imputation

- monotone missing data patterns

- non-monotone missing data patterns

- predictive mean matching

- regression methods

Multiple source data

Multi-stage sampling

Multivariate analysis of variance (MANOVA)

Multivariate normal distribution

Muscatine Coronary Risk Factor (MCRF) Study

National Cooperative Gallstone Study (NCGS)

National Health and Nutrition Examination Survey (NHANES)

Negative binomial distribution

Nested data

Nested models

New Haven Child Survey (NHCS)

*See also* Connecticut Child Surveys

Newton-Raphson algorithm

NIH method

*See also* Two-stage models

Non-linear regression

Normal distribution

Not missing at random (NMAR)

Nuisance parameters

Odds ratio

Offset

Onychomycosis Study

Ordinal regression models

adjacent-category model

continuation-ratio model  
illustration  
proportional odds model  
underlying latent variable  
Ordinary least squares (OLS)  
Outliers  
Overdispersion  
  extra-binomial variation  
  extra-Poisson variation

Parametric curves  
  general linear model formulation  
  linear trends  
  polynomial trends  
  quadratic trends

Pattern mixture model

Penalized quasi-likelihood  
  basis of approximation

Penalized splines  
  best linear unbiased predictor (BLUP)  
  mixed model representation  
  piecewise linear trend

Piecewise linear trend

Poisson distribution

Poisson regression  
  example  
  offset  
  overdispersion

Polynomial trends

Population-average models

Power  
  binary response  
  continuous response  
  impact of missing data

Prediction  
  shrinkage

Primary sampling unit (PSU)

Probit regression

Profile analysis  
  *See also* Analysis of response profiles

Promotion of Breastfeeding Intervention Trial (PROBIT)

Proportional hazards model

Proportional odds model

Quantile or Q-Q plot

Random effects  
  prediction  
  *See also* Best linear unbiased predictor

Randomized block design

- Rate ratio
- Reference group parameterization
- Reliability
- Repeated measurements
  - Repeated measures analysis by ANOVA
  - Repeated measures analysis by MANOVA
  - Repeated measures design
- Residual diagnostics
  - aggregating residuals
  - cumulative sum of residuals
  - transformed residuals
  - untransformed residuals
- Response profiles
- Response trajectories
- Restricted maximum likelihood (REML)
- Rotating panel design
- Sample size
  - binary response
  - continuous response
  - impact of missing data
- “Sandwich” variance estimator
- Satterthwaite approximation
- Scale parameter
- Semiparametric curves
  - linear splines
- Semiparametric regression models
- Semi-variogram
- Shrinkage
- Significance level
- Six Cities Study of Air Pollution and Health
- Skin Cancer Prevention Study
- Smoothing
  - bandwidth
  - bias and precision trade-off
  - penalized splines
  - roughness penalty
  - smoothing parameter
  - subject-specific smooth curves
  - truncated line function
- Spline
  - general linear model formulation
  - knot location
  - truncated line function
- Split-plot design
- Standard deviation
  - definition
- Standard error of measurement.
- Standardized mortality ratio (SMR)

Stationarity  
Statistical inference  
Study of Bladder Cancer Tumors  
Study of Dental Growth  
Study of Log C-Peptide Concentration in Children with Diabetes  
Study of Low Birth Weight Infants  
Study of Progesterone Metabolite Concentration  
Study of Risk Factors for Coronary Heart Disease (CHD)  
Study of Weight Gain  
Subject-specific models  
*See also* Generalized linear mixed effects models  
Substantive parameters  
Sufficient statistic  
Summary measure analysis  
Survival analysis  
Taylor series expansion  
Television, School and Family Smoking Prevention and Cessation Project  
Time plot  
Time series data  
Time-varying covariates  
    causal interpretation  
    endogenous  
    exogenous  
    external  
    fixed by design  
    interpretation of stochastic time-varying covariate effects  
    stochastic  
Toeplitz covariance  
Transformed residuals  
    Cholesky decomposition  
Treatment of Lead-Exposed Children (TLC) Trial  
Truncated line function  
Two-stage models  
Univariate repeated measures ANOVA  
Unstructured covariance  
U.S. Centers for Disease Control Study of Infant Birth Weight  
U.S. National Institute of Mental Health (NIMH) Schizophrenia Collaborative Study  
Variance  
    definition  
    heterogeneous over time  
    properties  
    residual diagnostics  
*See also* Overdispersion  
Variance function  
Vectors  
    arithmetic operations  
    basic concepts

definitions

transpose

Vlagtwedde-Vlaardingen Study

Wald test

multivariate Wald test

Within-individual biological variation

Within-individual change

Within-subject change

Within-subject variability

inherent biological variability

measurement error

Zero-inflated Poisson (ZIP) model