# ANALYSIS OF LONGITUDINAL AND CLUSTER-CORRELATED DATA

## Nan Laird
*Harvard University*

Conference Board of the Mathematical Sciences

*Regional Conference Series*
*in Probability and Statistics*

Supported by the
National Science Foundation

# Contents

# Preface

The analysis of data with outcomes measured repeatedly on each subject has experienced several transforming developments in the last twenty years. This monograph presents a unified treatment of modern methods for longitudinal and/or correlated data that have developed during this period. The basic approach that we take to modeling longitudinal data is to extend familiar univariate regression models to multivariate or correlated outcomes. We deal with linear models for measured data and generalized linear models for binary and count data. We show how methods can accommodate missing outcomes and/or unbalanced designs. Both likelihood and moment methods of estimation are covered, as are random effects approaches to data modeling and parameter estimation.

The monograph assumes that the reader has a solid foundation in statistical inference, linear and generalized linear regression models, and a basic knowledge of multivariate methods. It is appropriate for second year doctoral students or postdoctoral fellows in Statistics/Biostatistics as well as researchers or faculty interested in learning about the field.

# Acknowledgments

This monograph grew out of notes written for a course on longitudinal/multivariate data analysis taught by myself and others at the Harvard School of Public Health. I am grateful to many collegues and former students for their comments and contributions to the monograph. I am especially indebted to Andrea Rotnitzky, who wrote drafts of several sections of Chapters 1, 3 and 4. I am also grateful to Stuart Baker, Christl Donnelly, Garrett Fitzmaurice, Joe Hogan, Nick Lange, Stuart Lipsitz, Nick Horton and Jim Ware for collaborations over the years which have contributed greatly to my own work in this area. Roman Torgovitsky provided especially helpful comments on the monograph at a crucial point. Finally, the impetus to turn course notes into a monograph came when Paul Speckman organized a National Science Foundation workshop on longitudinal data at the University of Missouri at Columbia. I am grateful for the NSF support, and to NIGMS for supporting my research in longitudinal data during the last twenty years.