

Statistical Methods for Categorical Data Analysis Second Edition

Daniel A. Powers • Yu Xie

# STATISTICAL METHODS FOR CATEGORICAL DATA ANALYSIS: 2ND EDITION

# STATISTICAL METHODS FOR CATEGORICAL DATA ANALYSIS: 2ND EDITION

# **DANIEL A. POWERS**

Department of Sociology and Population Research Center, University of Texas at Austin, Austin, Texas, USA

# YU XIE

Department of Sociology, Department of Statistics, and Institute for Social Research, University of Michigan, Ann Arbor, Michigan, USA



United Kingdom • North America • Japan India • Malaysia • China Emerald Group Publishing Limited Howard House, Wagon Lane, Bingley BD16 1WA, UK

First edition 2008

Copyright © 2008 Emerald Group Publishing Limited

#### Reprints and permission service

Contact: booksandseries@emeraldinsight.com

No part of this book may be reproduced, stored in a retrieval system, transmitted in any form or by any means electronic, mechanical, photocopying, recording or otherwise without either the prior written permission of the publisher or a licence permitting restricted copying issued in the UK by The Copyright Licensing Agency and in the USA by The Copyright Clearance Center. No responsibility is accepted for the accuracy of information contained in the text, illustrations or advertisements. The opinions expressed in these chapters are not necessarily those of the Editor or the publisher.

#### British Library Cataloguing in Publication Data

A catalogue record for this book is available from the British Library

ISBN: 978-0-1237-2562-2



Awarded in recognition of Emerald's production department's adherence to quality systems and processes when preparing scholarly journals for print



To Our Parents

Dick and Janet Powers and Liangyao Xie and Huazhen Zhao

## Contents

Lis	ist of Figures	
Lis	List of Tables	
Pre	Preface	
1.	Introduction	1
2.	Review of Linear Regression Models	11
3.	Models for Binary Data	31
4.	Loglinear Models for Contingency Tables	67
5.	Multilevel Models for Binary Data	115
6.	Statistical Models for Event Occurrence	167
7.	Models for Ordinal Dependent Variables	221
8.	Models for Nominal Dependent Variables	243
Ap	ppendix A: The Matrix Approach to Regression	269
Ap	ppendix B: Maximum Likelihood Estimation	277
Re	References 2	
Inc	lex	307

# List of Figures

Figure 1.1:	Typology of the four types of measurements	5
Figure 2.1:	Maximization of log L with respect to $\theta$	20
Figure 2.2:	Logistic vs. linear regressions for binary data	29
Figure 3.1:	Logit and probit transformations of p	40
Figure 3.2:	Marginal effects as slopes of tangent lines to the cumulative probability curve	60
Figure 3.3:	Marginal effect of a dummy variable	62
Figure 3.4:	Graduation probability by family income	64
Figure 3.5:	Complementary log-log transformation of p	64
Figure 5.1:	Observed and predicted probabilities showing shrinkage of predicted probabilities toward the overall proportion	128
Figure 5.2:	Predicted probabilities of premarital birth by family structure and mother's education from model 2	132
Figure 5.3:	Distribution of family-specific random effects	134
Figure 5.4:	Empirical Bayes estimates of family-specific random effects	135
Figure 5.5:	Observed, marginal, and conditional logits	141
Figure 5.6:	Trace plots and histograms for $\beta_3$ and $\sigma_0^2$	148
Figure 5.7:	Posterior distribution of odds ratios of idleness (South versus non-South) (high-school graduate versus non-graduate)	153
Figure 5.8:	Item characteristic curves for 1PL with 3 items	159
Figure 5.9:	Item characteristic curve for 2PL	161

#### **x** List of Figures

Figure 5.10:	Item characteristic curve for 1PL and 2PL model using LSAT Data	164
Figure 6.1:	Discrete-time hazard and survivor functions for program dropout	174
Figure 6.2:	Plots of log cumulative hazard and survivor functions	211
Figure 6.3:	Plot of Schoenfeld residuals for the family income effect	214
Figure 6.4:	Plot of time varying family income effect	216
Figure 7.1:	Cumulative probabilities corresponding to a four-category response	229
Figure 7.2:	The relation between latent variables and realized outcomes	232

# List of Tables

Table 2.1:	Mortality in the first six months of life: Ümea, Sweden	25
Table 2.2:	Column-oriented layout of data file	25
Table 2.3:	OLS, FGLS, and ML estimates of the log-rate model	28
Table 2.4:	A typology of regression models	29
Table 3.1:	High school graduates by race, sex, and family structure	33
Table 3.2:	Column-formatted summary of data in Table 3.1 using dummy variables	34
Table 3.3:	Estimates from alternative binary response models	41
Table 3.4:	Estimated graduation probabilities by race, sex, and family structure	41
Table 3.5:	Comparing main effects and two-way interaction models	53
Table 3.6:	Voting inclination by income and sex	61
Table 3.7:	Logit and probit estimates from individual-level data	63
Table 4.1:	Education and attitude toward premarital sex	68
Table 4.2:	Observed (expected) frequencies	69
Table 4.3:	Expected probabilities	70
Table 4.4:	Expected frequencies under independence	71
Table 4.5:	Contribution to Pearson $\chi^2$	71
Table 4.6:	Row-specific proportions under independence	72
Table 4.7:	Row-specific proportions for observed data	72

### xii List of Tables

Table 4.8:	Full table for the attitude example	75
Table 4.9:	Local odds-ratios based on adjacent rows and columns	75
Table 4.10:	Pearson $\chi^2$ components under model A	78
Table 4.11:	Identifiable parameters	84
Table 4.12:	Hauser's mobility table	87
Table 4.13:	Interaction parameters of the saturated model: intergenerational mobility example	88
Table 4.14:	Estimated $\mu^h$ parameters	89
Table 4.15:	Goodness-of-fit statistics for mobility models	97
Table 4.16:	Attitudes toward abortion and premarital sex	100
Table 4.17:	Estimated scale scores	101
Table 4.18:	Graduate admission data from UC-Berkeley	102
Table 4.19:	Collapsed graduate admission data	103
Table 4.20:	Goodness-of-fit statistics of models for admission data	108
Table 4.21:	Estimates of interaction parameters of model 4	108
Table 4.22:	Models for three-nation class mobility data	111
Table 4.23:	Nation-specific $\phi$ parameters	113
Table 5.1:	Post-doctorate training in biochemistry and NIH funding	124
Table 5.2:	Conventional and random intercept models	125
Table 5.3:	Observed $(\tilde{p})$ and model-based $(\hat{p})$ proportions	127
Table 5.4:	Logit models for first premarital birth	131
Table 5.5:	Points $(u)$ and weights $(p)$ for numerical integration	133
Table 5.6:	Longitudinal models fit to Youth Employment Data	140
Table 5.7:	Estimates from different approaches	149
Table 5.8:	Observed and expected response patterns	152

Table 5.9:	Fit statistics for logit models	154
Table 5.10:	Bock and Lieberman Law School Aptitude Test (LSAT) Data	163
Table 5.11:	1PL and 2PL models estimated from LSAT Data	164
Table 6.1:	Event occurrence data	171
Table 6.2:	Life table for program dropout	173
Table 6.3:	Person-level and person-period format data	178
Table 6.4:	Sequences of observed binary responses over a five-wave panel	178
Table 6.5:	Estimates from discrete-time logit model of program dropout	182
Table 6.6:	Waiting time to program completion	185
Table 6.7:	Occurrence-exposure matrix for data in Table 6.6	187
Table 6.8:	Infant mortality (exposure in days) by age, race, and birth outcomes in the U.S. — 1995–1998	191
Table 6.9:	Models and fit statistics for infant mortality data	193
Table 6.10:	Baseline hazards and hazard ratios of infant mortality in the U.S. — 1995–1998	194
Table 6.11:	Conceptual format for event-history data	195
Table 6.12:	Arrangement of split-episode event-history data	198
Table 6.13:	Cross-classified data from Table 6.12	199
Table 6.14:	Piecewise constant rate model with nonproportional effects	199
Table 6.15:	Piecewise constant exponential models for risk of first premarital birth	200
Table 6.16:	Diagnostic tests for nonproportionality	214
Table 6.17:	Cox regression models with proportional and nonproportional effects	215
Table 7.1:	Normal score transformation for the attitude example	224
Table 7.2:	Education and attitude toward premarital sex	226
Table 7.3:	Ordered logit estimates under alternative parameterizations	235

#### xiv List of Tables

Table 7.4:	Ordered probit estimates and marginal effects	237
Table 7.5:	Attitude towards women's employment	238
Table 7.6:	Ordered logit and separate logit estimates	238
Table 7.7:	Brant tests of proportional odds assumption	239
Table 7.8:	Partial proportional odds model	240
Table 8.1:	Multinomial logit results	250
Table 8.2:	Equivalence between multinomial logit and loglinear models for three-way tables	253
Table 8.3:	Employment status by race and father's education	254
Table 8.4:	Multinomial logit estimates derived from a loglinear model	254
Table 8.5:	Estimates from conditional logit model	258
Table 8.6:	Results from the mixed model	260
Table 8.7:	Educational attainment	267

## Preface

In this book, we intend to give a comprehensive introduction to methods and models for the analysis of categorical data and their applications in social science research. The primary audiences are graduate students and practicing researchers in social science. The book also serves as a reference.

One feature that distinguishes our book from other books on the topic is our explicit aim to integrate the transformational approach and the latent variable approach, two diverse but complementary traditions dealing with the analysis of categorical data. The statistical, or transformational, approach to categorical data analysis is most familiar to researchers in demography and biostatistics, whereas the latent variable approach is often taken by economists. A discussion of the two approaches is given in Chapter 1.

We assume that the reader has prior knowledge such as that covered in a typical applied regression course but not necessarily in advanced mathematical statistics. Although some technical details are unavoidable in a book like this, we make the book accessible by resorting to substantive examples. Some readers may wish to skip portions of the book that are technical without losing much appreciation of the book.

To utilize the internet technology fully, we have set up a website for the book at http://webspace.utexas.edu/dpowers/www/.<sup>1</sup>

The website contains the data sets and the programming codes for the examples discussed in the book in several statistical packages including: GLIM (Numerical Algorithms Group, 1986), LIMDEP (Greene, 2007), SAS (SAS Institute, 2004), Stata (Stata Corporation, 2007), TDA (Rohwer & Pötter, 2000), and R (R Development Core Team, 2006). The website provides some GLIM macros and GAUSS (Aptech Systems, 1997) and R subroutines to illustrate details of estimation as well as several applications of specialized programs for models that cannot be estimated in a standard software package, for example aML (Lillard & Panis, 2003). We will continue to update the website as new programs become available.

<sup>1.</sup> This page is linked at YuXie.com and Powers-Xie.com.

### New to the 2nd Edition

We have updated the material in each chapter and have included a new chapter on multilevel models for binary data (Chapter 5). This chapter provides details on marginal maximum likelihood estimation and modern Bayesian estimation methods. We include a discussion of Rasch models and random-coefficient models for longitudinal analysis. We have reorganized the chapter on event-history models (Chapter 6) and include expanded coverage of discrete-time models and Cox regression models. The chapters on ordinal (Chapter 7) and nominal (Chapter 8) response models have also been updated.

### Use of This Text in a Course on Categorical Data Models

This book is appropriate for a single-term course in categorical data modeling. Chapters 1 and 2 provide an introduction and basic foundation for the course. Our view is that, regardless of the type of data, a regression-type modeling approach can be an appropriate analytic method. Chapter 3 provides an introduction and detailed treatment of regression models for binary data. Chapter 4 goes into greater detail on the methods for analyzing contingency tables. Chapter 5 discusses multilevel/ hierarchical models for binary data. Chapter 6 covers event history techniques. Chapters 7 and 8 provide an overview of methods for ordered and unordered categorical response variables. This material is linked to the contingency table approach of Chapter 4 and the latent variable framework outlined in Chapter 3.

#### Acknowledgments

At various stages of the book project, we benefited from the encouragement of and association with the following scholars: Paul Allison, Mark Becker, John Fox, Richard Gonzalez, Leo Goodman, David Grusky, Robert Hauser, Michael Hout, Kenneth Land, Scott Long, Charles Manski, Robert Mare, Bill Mason, Susan Murphy, Trond Peterson, Thomas Pullum, Adrian Raftery, Steve Raudenbush, Arthur Sakamoto, Herbert Smith, Michael Sobel, Chris Winship, Raymond Wong, Larry Wu, and Kazuo Yamaguchi. In addition, we extend our gratitude to many graduate students who have taken statistics courses from us and inspired us to write the book.

A Dean's Fellowship at the University of Texas at Austin to Dan Powers, a National Science Foundation's Young Investigator Award to Yu Xie, and University of Michigan's internal funds to Yu Xie provided partial support for this project.

We also thank the external reviewers for providing valuable critiques on early versions of the manuscript and Pam Bennett, John Fox, Kimberly Goyette, and James Raymo for carefully proofreading the final version of the manuscript and providing programming examples in the first edition. We thank Meichu D. Chen and the many graduate students who spotted errors in the first edition. Special thanks go to Cathy (Hui) Liu who carefully read the new material for the 2nd edition. We also appreciate the excellent editorial assistance of Cindy Glovinsky. We alone are responsible for errors that remain.

Last, but not least, we thank our editor J. Scott Bentley at Academic Press and Elsevier for initiating the project and then making efforts to bring the first edition to completion, as well as for guiding the development of the 2nd edition. We thank Rachel Brown at EmeraldInsight for her help in producing the 2nd edition, as well as Macmillan for assistance in formatting the book.

> Daniel A. Powers Yu Xie

### Chapter 1

## Introduction

#### 1.1. Why Categorical Data Analysis?

What is common about birth, marriage, schooling, employment, occupation, migration, divorce, and death? The answer: they are all categorical variables commonly studied in social science research. In fact, most observed outcomes in social science research are measured categorically. If you are a practicing social scientist, chances are good that you have studied a phenomenon involving a categorical variable. (This is true even if you have not used any special statistical method for handling categorical data.) If you are in a graduate program to become a social scientist, you will soon, if not already, encounter a categorical variable. Notice that even our statement of whether or not you have encountered a categorical variable in your career is itself a categorical measurement!

Statistical methods and techniques for categorical data analysis have undergone rapid development in the past 25 years or so. Their applications in applied research have become commonplace in recent years, due in large part to the availability of commercial software and inexpensive computing. Since some of the material is rather new and dispersed among several disciplines, we believe that there is a need for a systematic treatment of the subject in a single book. This book is aimed at helping applied social scientists use special tools that are well suited for analyzing categorical data. In this chapter, we will first define categorical variables and then introduce our approach to the subject.

#### 1.1.1. Defining Categorical Variables

We define categorical variables as those variables that can be measured using only a limited number of values or categories. This definition distinguishes categorical variables from continuous variables, which, in principle, can assume an infinite number of values.

Although this definition of categorical variables is clear, its application to applied work is far more ambiguous. Many variables of long-lasting interest to social scientists are clearly categorical. Such variables include: race, gender, immigration status, marital status, employment, birth, and death. However, conceptually continuous variables are sometimes treated as continuous and other times as categorical. When a continuous variable is treated as a categorical variable, it is called *categorization* or *discretization* of the continuous variable. Categorization is often necessary in practice because either the substantive meaning or the actual measurement of a continuous variable is categorical. Age is a good example. Although conceptually continuous, age is often treated as categorical in actual research for substantive and practical reasons. Substantively, age serves as a proxy for qualitative states for some research purposes, qualitatively transforming an individual's status at certain key points. Changes in legal and social status occur first during the transition into adulthood and later during the transition out of the labor force. For practical reasons, age is usually reported in single-year or five-year intervals.<sup>1</sup>

Indeed, our usual instruments in social science research are crude in the sense that they typically constrain possible responses to a limited number of possible values. It is for this reason that we earlier stated that most, if not all, observed outcomes in social science are categorical.

What variables should then be considered categorical as opposed to continuous in empirical research? The answer depends on many factors, two of which are their substantive meaning in the theoretical model and their measurement precision. One requirement for treating a variable as categorical is that its values are repeated for at least a significant portion of the sample.<sup>2</sup> As will be shown later, the distinction between continuous and categorical variables is far more consequential for response variables than for explanatory variables.

#### 1.1.2. Dependent and Independent Variables

A *dependent* (also called response, outcome, or endogenous) variable represents a population characteristic of interest being explained in a study. *Independent* (also called explanatory, predetermined, or exogenous) variables are variables that are used to explain the variation in the dependent variable. Typically, the characteristic of interest is the population mean of the dependent variable (or its transformation) *conditional* on values of an independent variable or set of independent variables. It is in this sense that we mean that the dependent variable depends on, is explained by, or is a function of independent variables in regression-type statistical models.

By *regression-type statistical models*, we mean models that predict either the expected value of the dependent variable or some other characteristic of the dependent variable, as a regression function of independent variables. Although in principle we could design our models to best predict any population parameter (e.g., the median) of the dependent variable or its transformation, in practice we

<sup>1.</sup> Education is another example. The substantive distinctions among "less than 12 years of schooling," "high-school diploma," "college degree," or "graduate degree" cannot be captured without categorization.

A few categories offer a concise representation of the important points in the distribution of education. 2. Note that a continuous variable can be truncated, meaning that it has zero probability of yielding a value beyond a particular threshold or cut-off point. When a continuous variable is truncated, the untruncated part is still continuous, whereas the part that is truncated resembles a categorical variable.

commonly use the term *regression* to denote the problem of predicting conditional means. When the regression function is a linear combination of independent variables, we have so-called linear regressions, which are widely used for continuous dependent variables.

#### 1.1.3. Categorical Dependent Variables

Although categorical and continuous variables share many properties in common, we wish to highlight some of the differences here. The distinction between categorical and continuous variables as dependent variables requires special attention. In contrast, the distinction is of relatively minor significance when they are used as independent variables in regression-type statistical models. Our definition of regression-type statistical models includes statistical methods for the analysis of variance and covariance, which can be represented by regressing the dependent variable on a set of dummy variables and, in the case of the analysis of covariance, other continuous covariates. Hence, including categorical variables as independent variables in regression-type models does not present any particular difficulties, as it mainly involves constructing dummy variables corresponding to different categories of the independent variable; all known properties of regression models are directly generalizable to models for the analysis of variance and covariance. As we will show later in this book, the situation changes drastically when we treat categorical variables as dependent variables, as much of our knowledge derived from linear regressions is simply inapplicable. In brief, special statistical methods are required for categorical data analysis (i.e., analysis involving categorical dependent variables).

Although the methods for analyzing categorical variables as independent variables in regression-type models have been a part of the standard statistical knowledge base that is now required for most advanced degrees in social science, methods for the analysis of categorical dependent variables are much less widely known. Much of the fundamental research on the methodology of analyzing categorical data has been developed only recently. We aim to give a systematic treatment of several important topics on categorical data analysis in this book so as to facilitate the integration of the material into social science research.

Unlike methods for continuous variables, methods for categorical data require close attention to the type of measurement of the dependent variable. Methods for analyzing one type of categorical dependent variable may be inappropriate for analyzing another type of variable.

#### 1.1.4. Types of Measurement

The type of measurement plays a key role in determining the appropriate method of analysis when a variable is used as a dependent variable. We present a typology for four types of measurement based on three distinctions.<sup>3</sup> First, let us distinguish between *quantitative* and *qualitative* measurements. The distinction between the two is that quantitative measurements closely index the substantive meanings of a variable with numerical values, whereas numerical values for qualitative measurements are substantively less meaningful, sometimes merely as classifications to denote mutually exclusive categories of characteristics (or attributes) uniquely. Qualitative variables are categorical variables.

Within the class of *quantitative* variables, it is often useful to distinguish further between *continuous* and *discrete* variables. Continuous variables, also called interval variables, may assume any real value. Variables such as income and socioeconomic status are typically treated as continuous over their plausible range of values. Discrete variables may assume only integer values and often represent event counts. Variables such as the number of children per family, the number of delinquent acts committed by a juvenile, and the number of accidents per year at a particular intersection are examples of discrete variables. According to our earlier definition, discrete (but quantitative) variables are also categorical variables.

Qualitative measurements can be further distinguished between ordinal and nominal. Ordinal measurements give rise to ordered qualitative variables, or ordinal variables. It is quite common to use numerical values to denote the ordering information in an ordered qualitative variable. However, numerical values corresponding to categories of ordinal variables reflect only the ranking order in a particular attribute; therefore, distances between two adjacent values are not the same. Attitudes toward gun control (strongly approve, approve, neutral, disapprove, and strongly disapprove), occupational skill level (highly skilled, medium skilled, low skilled, and unskilled), and the classification of levels of education as (grade school, high school, college, and graduate) are examples of ordinal variables.

Nominal measurements yield unordered qualitative variables, often referred to as *nominal* variables. Nominal variables possess no inherent ordering, nor numerical distance, between category levels. Classifications of race and ethnicity (white, black, Hispanic, and other), gender (male and female), and marital status (never married, married, divorced, and widowed) are examples of unordered qualitative variables. It is worth noting at this point, however, that the distinction between ordinal and nominal variables is not always clear-cut. Much of the distinction depends on the research questions. The same variable may be ordinal for some researchers but nominal for others.

To further illustrate the last point, let us use occupation as an example. Distinct occupations are often measured by open-ended questions and then manually coded into a classification system with three-digit numerical codes that do not represent magnitudes in substantive dimensions. Since the number of potential occupations is large (usually at least a few hundred in a coding scheme for a modern society), it is desirable, and indeed necessary, to reduce the amount of detail in an occupational

<sup>3.</sup> For an historical background, see Duncan's (1984) important book Notes on Social Measurement.



Figure 1.1: Typology of the four types of measurements.

measure through data reduction. One method of data reduction is to collapse detailed occupational codes into major occupational categories and treat them as constituting either an ordinal or even a nominal measurement (Duncan, 1979; Hauser, 1978). Another method of data reduction is to scale occupations along the dimension of a socioeconomic index (SEI) (Duncan, 1961) — thus into an interval variable. More recently, Hauser and Warren (1997) challenged Duncan's approach and suggested instead that to measure occupational socioeconomic status, occupational education. Hauser and Warren's work illustrates the importance of considering multiple dimensions when nominal measures are scaled into interval measures.

Figure 1.1 summarizes our typology scheme for the four types of measurements. According to this typology, there are three types of categorical variables: discrete, ordinal, and nominal, all of which will be discussed in this book. This distinction among the three types of categorical variables is useful only when the number of possible values equals or exceeds three. When the number of possible values is two, we have a special case called a binary variable. A *binary* variable can be discrete, ordinal, or nominal, depending on the researcher's interpretation. For example, if a researcher is interested in studying compliance with the one-child policy in China, the dependent variable is whether a couple has given birth to more than one child. For simplicity, assume that in a particular sample a woman has at least one child and no more than two children. Let us code y so that y = 0 if a woman has one child, and y = 1 if she has two children. In this case, the dependent variable can be interpreted as discrete (number of children -1), ordinal (one child or more than one child), or nominal (compliance vs. noncompliance). Fortunately, the researcher may apply the same statistical methods for all three cases. It is the substantive understanding of the results that varies from one interpretation to another.

#### 1.2. Two Philosophies of Categorical Data

The development of methods for the analysis of categorical data has benefitted greatly from contributions by scholars in such diverse fields as statistics, biostatistics, economics, psychology, and sociology. This multidisciplinary origin has given categorical data analysis multiple approaches to similar problems and multiple interpretations for similar methodologies. As a result, categorical data analysis is an intellectually rich and expanding field. However, this interdisciplinary nature has also

made synthesizing and consolidating available techniques difficult due to the diverse applications and differing terminology across disciplines.

Part of this difficulty stems from two fundamentally different "philosophies" concerning the nature of categorical data. One philosophy views categorical variables as being inherently categorical and relies on transformations of the data to derive regression-type models. The other philosophy presumes that categorical variables are conceptually continuous but are observed, or measured, as categorical. In the onechild policy example, a researcher may view "compliance" as a behavioral continuum. However, he/she can only observe two distinct values of this dependent variable. This approach relies on latent variables to derive regression-type models. These very different philosophies can be traced back to the acrimonious debate between Karl Pearson and G. Udny Yule between 1904 and 1913 (Agresti, 2002, pp. 619–622). Although these two approaches can be found in any single discipline, the first is more closely identified with statistics and biostatistics, and the second with econometrics and psychometrics. For simplicity, we will refer to the first approach as statistical or transformational and to the second as econometric or latent variable. We intend the terms statistical and econometric here as short-hand labels rather than as descriptions of the two disciplines.

#### 1.2.1. The Transformational Approach

In the *transformational*, or statistical, approach, categorical data are considered as inherently categorical and should be modeled as such. In this approach, there is a direct one-to-one correspondence between population parameters of interest and sample statistics. The focus is on estimating population parameters that correspond to their sample analogs. No latent, or unobserved, variable is invoked.

In the transformational approach, statistical modeling means that the expected value of the categorical dependent variable, after some transformation, is expressed as a linear function of the independent variables. Given the categorical nature of the dependent variable, the regression function cannot be linear. The problem of nonlinearity is handled through nonlinear functions that transform the expected value of the categorical variable into a linear function of the independent variables. Such transformation functions are now commonly referred to as *link* functions.<sup>4</sup>

For example, in the analysis of discrete (count) data, the expected frequencies (or cell counts) must be nonnegative. To ensure that the predicted values from regression models fit these constraints, the natural logarithm function (or *log* link) is used to transform the expected value of the dependent variable so that a model for the logged count can be expressed as a linear function of independent variables. This *loglinear* transformation serves two purposes: it ensures that the fitted values are appropriate

<sup>4.</sup> Models that can be transformed to linear models via link functions are referred to as *generalized linear models*. McCullagh and Nelder (1989) provide an extensive treatment of these types of models.

for count data (i.e., nonnegative), and it permits the unknown regression parameters to lie within the entire real space (parameter space).

In binomial response models, estimated probabilities must lie in the interval [0,1], a range that is violated by any linear function if independent variables are allowed to vary freely. Instead of directly modeling probabilities in this range, we can model a transformation of probability that lies in the interval  $(-\infty, +\infty)$ . There are a number of ways to transform probabilities. The *logit* transformation,  $\log [p/(1-p)]$ , can be used to transform the probability scale so that it can be expressed as a linear function of independent variables. A *probit* transformation,  $\Phi^{-1}(p)$ , can be used in a similar fashion to re-scale probabilities. The probit link utilizes the inverse of the cumulative standard normal distribution function to transform the expected probability to the range  $(-\infty, +\infty)$  (i.e., by transforming probabilities to Z-scores). As in the logit model, the probit link transforms the probability so that it can be expressed as a linear function of independent variables. Both the logit and probit transformations ensure that the predicted probabilities are in the proper range for all possible values of parameters and independent variables.

#### 1.2.2. The Latent Variable Approach

The latent variable, or econometric, approach provides a somewhat different view of categorical data. The key to this approach is to assume the existence of a continuous unobserved or *latent* variable underlying an observed categorical variable. When the latent variable crosses a threshold, the observed categorical variable takes on a different value. According to the latent variable approach, what makes categorical variables different from usual continuously distributed variables is partial observability. That is, we can infer from observed categorical values only the intervals within which latent variables lie but not the actual values themselves. For this reason, econometricians commonly refer to categorical variables as limited-dependent variables (Maddala, 1983).

In the latent variable approach, the researcher's theoretical interest lies more in how independent variables affect the latent continuous variables (called structural analysis) than in how independent variables affect the observed categorical variable. From the latent variable perspective, it is thus convenient to think of the sample data as actual *realizations* of population quantities that are *unobservable*. For instance, the observed response categories may reflect the actual choices made by individuals in a sample, but underlying each choice at the population level is a latent variable representing the difference between the cost and the benefit of a particular choice made by an individual decision maker. Similarly, a binary variable may be thought of as the sample realization of a continuous variable representing an unobserved *propensity*. For example, in studies of college admissions, we may assume the existence of a continuous latent variable — qualification — such that applicants whose qualifications exceed the required threshold are admitted, and those whose qualifications fall short of the threshold are rejected (Manski & Wise, 1983).

In studies of women's labor force participation, economic reasoning holds that a woman will participate in the labor force if her market wage exceeds her reservation wage (Heckman, 1979). In practice, it is not possible for the researcher to observe applicants' qualifications, nor the difference between the market and reservation wages. We can, however, observe admission decisions and labor force participation status, which can be taken as *observed* realizations of the underlying population-level latent variable representing likelihood of admission or labor force participation.

Experimental studies in the biological sciences have also made good use of latent variables. In studies of the effectiveness of pesticides, for example, whether an insect dies depends on its *tolerance* to a level of dosage of an insecticide. It is assumed that an insect will die if a dosage level exceeds the insect's tolerance. The binary variable (lives/dies) is the realization of a continuous unobservable variable, the difference between dosage and tolerance.

The latent variable concept has been extended to the construction of latent *categorical* variables. A prime example is the latent class model, which capitalizes on independence conditional on membership in latent classes. This is analogous to factor analysis for continuously distributed variables. Heckman and Singer's (1984) nonparametric method of handling unobserved heterogeneity in survival analysis is also rooted in this fundamental idea.

#### **1.3.** An Historical Note

The development of techniques for the analysis of categorical data has been motivated in part by particular substantive concerns in fields such as sociology, economics, epidemiology, and demography (for an historical account in social science, see Camic & Xie, 1994). For example, several innovations in loglinear modeling had their origins in the study of social mobility (e.g., Duncan, 1979; Goodman, 1979; Hauser, 1978); the literature on sample selection models emerged from economic analyses of women's earnings (Heckman, 1979); and problems in the analysis of consumer choices led to the development of many of the techniques for multicategory response variables (McFadden, 1974). Methodological advances in survival analysis arose as extensions of the life-tables technique in demography by statisticians and biostatisticians to incorporate covariates in modeling hazard rates (Cox, 1972; Laird & Oliver, 1981). McCullagh and Nelder's (1989) theory of generalized linear models provided a unified framework which can be applied to most of these models.

Today's latent variable approach grew out of the early psychophysics tradition, where observed frequency distributions of qualitative "judgments" were used to scale the intensity of continuously distributed stimuli (e.g., Thurstone, 1927). In the experimental framework of psychophysics, the "latent" variables were unobservable only to the subjects under an experiment, since the stimuli were manipulated by and thus known to the researcher. For illustration, imagine that a group of subjects are asked to rank the relative weights of two similar objects given by the experimenter. It is reasonable to assume that the probability of giving the correct answer is positively associated with the actual difference in weight. Thurstone (1927) explicitly assumed a normal distribution for the psychological stimulus and related it to the distribution of "judgments," thus paving the way to today's probit analysis. With time, social scientists have expanded this approach to uncover properties of latent variables from observed data, through such techniques as latent trait models and latent class models. For a treatment of sociologists' contributions to the latent variable approach, see Clogg (1992).

### 1.4. Approach of This Book

Two features distinguish this book from other texts on the analysis of categorical data. First, this book presents both the transformational and latent variable approaches and, in doing so, synthesizes similar methods in statistical and econometric literatures. Whenever possible, we shall show how the two approaches are similar and in what ways they are different. Second, this book has an applied as opposed to theoretical orientation. We shall draw examples from applied social science research and use data sets constructed for pedagogical purposes. In keeping with the applied orientation of this book, we shall also present actual programming examples for the models discussed, while keeping theoretical discussions at a minimum. We shall provide our data sets, program code, and computer outputs through a website.<sup>5</sup>

#### 1.4.1. Combining the Statistical and Latent Variable Approaches

In many instances, the transformational and latent variable approaches are simply two parallel ways of looking at the same phenomena. More often than not, the two approaches yield exactly the same statistical procedures except for minor differences due to the manner in which the model is specified or parameterized. When this is the case, one's viewpoint about the underlying nature of observed categorical variables does not affect specific statistical techniques that we will cover but simply alters the substantive interpretations of results.

### 1.4.2. Organization of the Book

This book begins by considering the simplest models for categorical data and proceeds to more complex models and methods. We begin with a review of the

<sup>5.</sup> Our website is continuously updated with new examples utilizing several computer packages. The URL is http://webspace.utexas.edu/dpowers/www/, linkable through YuXie.com and Powers-Xie.com.

general concepts behind regression models for continuous dependent variables. This is a natural starting point since many of the familiar ideas and principles used in the analysis of covariance and regression for continuous variables will carry over to the analysis of categorical dependent variables. These concepts are described in Chapter 2, along with a general orientation to regression models. Chapter 3 discusses models for binary data and issues pertaining to estimation, model building, and the interpretation of results. Chapter 4 provides an overview of measures of association and models for contingency tables. Chapter 5 builds on material in Chapter 3 to introduce multilevel (or hierarchical) models for binary data. Chapter 6 presents methods for event occurrences in time. Chapters 7 and 8 outline various methods for the analysis of polytomous (or multinomial) response variables that assume ordinal or nominal measures.

# **Review of Linear Regression Models**

#### 2.1. Regression Models

This chapter reviews the classic linear regression model for continuous dependent variables. We assume the reader's familiarity with the linear regression model and thus will not delve into its details. Instead, we will highlight some general concepts and principles underlying the linear regression model that will be useful in later chapters focused on categorical dependent variables.

Regression is one of the most widely used statistical techniques for analyzing observational data. As mentioned in Chapter 1, the analysis of observational data typically requires a structural and multivariate approach. Regression models are used in this context to uncover net relationships between an outcome, or response, variable and a few key explanatory variables while controlling for confounding factors.

Regression models are used to meet different research goals. Sometimes, regression modeling is aimed at learning the causal effect of one variable, or a set of variables, on a dependent variable. Other times, regression models are used to predict the value of a response variable. Finally, regression models are often intended as short-hand summaries providing a description linking a dependent variable and independent variables.

#### 2.1.1. Three Conceptualizations of Regression

A researcher faced with a large amount of raw data will want to summarize it in a way that presents essential information without too much distortion. Examples of data reduction include frequency tables or group-specific means and variances. Like most methods in statistics, regression is also a data-reduction technique. In regression analysis, the objective is to predict, as closely as possible, an array of observed values of the dependent variable based on a simple function of independent variables. Obviously, predicted values from regression models are not exactly the same as observed ones. Characteristically, regression partitions an observation into two parts:

The observed part represents the actual values of the dependent variable at hand. The structural part denotes the relationship between the dependent and independent variables. The stochastic part is the random component unexplained by the structural part. In general, the last term may be regarded as the sum of three components: omitted structural factors, measurement error, and "noise." Omitting structural factors is inevitable in social science research because we can never claim to understand and measure all causal structures affecting a dependent variable. Measurement error refers to inaccuracies in the way in which the data are recorded, reported, or measured. Random noise reflects the extent to which human behavior or occurrence of events is subject to uncertainty (i.e., stochastic influences).

How to interpret regression models is contingent on one's conceptualization about what regression does to data. We propose three different conceptualizations.

Causation : observed = true mechanism + disturbance
Prediction : observed = predicted + error
Description : <u>observed</u> = <u>summary</u> + <u>residual</u>

These conceptualizations provide three different views of quantitative analysis. The first approach corresponds most closely to what might be perceived as a view in classical econometrics in which the model accurately represents the "true" causal mechanism that generates the data. The researcher's goal is to specify a model to uncover the data-generating mechanism, or "true" causal model. This first approach can be viewed as an attempt to get as close as possible to a deterministic model. More modern approaches would argue that there is no "true" model but rather that some models are more useful, more interesting, or closer to the truth than others.

The second approach is more directly applicable to fields like engineering where, given a relationship between explanatory variables and a response variable, the goal is to make useful response predictions for new data. For example, suppose that the strength of a material is related to temperature and pressure during the manufacturing process. Suppose that we produce a sample of materials by varying temperature and pressure in a systematic way. One objective of modeling might be to find the values of temperature and pressure that give the material maximum strength. Social scientists also employ this modeling approach in forecasting and may use this approach to identify people at risk of a particular outcome based on certain characteristics.

The third approach reflects the current view in modern econometrics and statistics in which a model serves to summarize the basic features of data without distorting them. A principle called Occam's razor, or the *law of parsimony*, is often invoked when assessing competing explanations of the same phenomenon. When applied to statistical models, this principle means that if two models equally explain the observed facts, the simpler model is preferred until new evidence proves otherwise. This approach differs from the first view in the sense that the question asked is not whether the model is "true" but whether it corresponds to the facts. The facts usually require formalization based on past research or theory. The model is then specified in accordance with theory or previous research. These conceptualizations are not mutually exclusive; the applicability of a particular interpretation hinges on concrete situations, particularly the nature of the research design and objectives. With most applications in social sciences utilizing observational data, our inclination is to favor the last interpretation (Xie, 2007). That is, the primary goal of statistical modeling is to summarize massive amounts of data with simple structures and few parameters. With this conceptualization of regression models, it is important to keep in mind the trade-off between accuracy and parsimony. On the one hand, we desire accuracy in a model in the sense that we want to preserve maximum information and minimize errors associated with residuals. On the other hand, we prefer parsimonious models. More often than not, the desire to preserve information can only be achieved by building complicated models, which comes at the expense of parsimony or simplicity. The tension between accuracy and parsimony is so fundamental to social science research that we will revisit the issue several times in the book.

#### 2.1.2. Anatomy of Linear Regression

There are three types of variables in a regression model: a dependent variable, a set of independent variables, and random errors. Because the exact nature of the dependency of the dependent variable on the independent variables is unknown, researchers often summarize it as a linear relationship in an approximation involving a set of unknown parameters or coefficients.

The continuous dependent variable, also called the response variable, is usually denoted by y. For a given sample of size n, we denote the individual data values as  $y_i = y_1, y_2, ..., y_n$ . We can think of the many possible values of y as forming a *population*. Like all random variables, y has a mean, a variance, and additional parameters to describe its distribution. The mean or expected value of y is denoted by  $E(y) = \mu$ . We can also let the mean of y be expressed as a function of independent variables. For example, if an independent variable assumes a unique value for each element in the population, and E(y) is modeled as a function of the independent variable, there would be a different mean, say  $\mu_i$ , for each observation.

More generally, associated with each observation is a set of independent variables, also called explanatory variables. The set of independent variables constitutes a data matrix indexed by *n* rows — corresponding to *n* individual units of analysis — and K + 1 columns — corresponding to *K* distinct independent variables plus a constant.<sup>1</sup> We will denote the  $n \times (K + 1)$  matrix of independent variables as **X**, where *K* is the total number of explanatory variables. The values of **X** for the *i*th observation are denoted by the vector  $\mathbf{x}_i = (x_{i0}, x_{i1}, \dots, x_{iK})^t$ . With no loss of

<sup>1.</sup> Throughout this book we will use bold-faced symbols to indicate that a quantity is a matrix or vector. When possible, we will use the more familiar "scalar" representations. Some basic principles of matrix algebra are reviewed in Appendix A.

generality, we include as the first column of **X** a vector of 1's (e.g.,  $x_{i0} = 1$ ), whose coefficient is the intercept. We may write the expression for the mean of *y*, conditional on the independent variables as

$$E(y_i | \mathbf{x}_i) = \beta_0 + \beta_1 x_{i1} + \beta_2 x_{i2} + \dots + \beta_K x_{iK}$$
$$= \sum_{k=0}^K \beta_k x_{ik}$$
$$= \mathbf{x}'_i \boldsymbol{\beta}$$
(2.1)

The  $\beta_k$  (k = 0, ..., K) terms are unknown regression coefficients, or parameters, to be estimated from the sampled data. The intercept,  $\beta_0$ , can be interpreted as the mean of y when all x variables are zero. The remaining  $\beta_k$  (k = 1, ..., K) terms are regression slopes, reflecting the amount that E(y) changes when  $x_{ik}$  changes by one unit, while holding other independent variables constant. The symbol  $\beta$  is used to denote the  $(K + 1) \times 1$  vector of regression coefficients,  $\beta = (\beta_0, \beta_1, ..., \beta_K)'$ .

In focusing on the expected value of y, other characteristics of the distribution of y are usually ignored. Since a model based on a set of independent variables cannot predict exactly the observed values of y, it is necessary to introduce  $\varepsilon_i$  (i.e., error, disturbance, or residual, depending on one's viewpoint). For the *i*th observation, we have

$$y_{i} = \beta_{0} + \beta_{1}x_{i1} + \beta_{2}x_{i2} + \ldots + \beta_{K}x_{iK} + \varepsilon_{i}$$
$$= \sum_{k=0}^{K} \beta_{k}x_{ik} + \varepsilon_{i}$$
$$= \mathbf{x}_{i}'\boldsymbol{\beta} + \varepsilon_{i}$$
(2.2)

This expression describes the way in which y is decomposed into a linear function of x's with unknown parameters ( $\beta$ ) and a residual term ( $\varepsilon_i$ ). Since  $\varepsilon_i$  is intrinsically unobservable, simplifying assumptions about the characteristics of  $\varepsilon_i$  are necessary. A key assumption that yields the identification of the unknown parameters in Equation 2.2 is the independence between  $\varepsilon$  and the x variables. Other assumptions are often invoked to improve efficiency. For example, it is common to assume  $\varepsilon_i$  to be independent of one another and identically distributed (i.i.d.). The independence assumption implies that the correlation in  $\varepsilon$  between a pair of observations is zero, whereas the identical distribution assumption assures a common variance of  $\sigma_{\varepsilon}^2$ (i.e., homoscedasticity). With the i.i.d. assumption, Eq. 2.2 can be estimated using ordinary least squares (OLS), which is described in Section 2.2.1.

Even without the i.i.d. assumption, however, the OLS estimator is still a consistent estimator if  $\varepsilon$  is uncorrelated with the x's, meaning that it converges to the parameter vector when the sample size is large.

#### 2.1.3. Basics of Statistical Inference

To understand estimation and statistical inference, it is necessary to introduce the distinction between population quantities (*parameters*) and their sample counterparts (*statistics*). This distinction is the basis for statistical inference, the practice of inferring characteristics of a population from more limited information contained in a sample drawn from the population. We begin with a general discussion of inference, although in this book inference is more narrowly limited to the estimation of parameters and their standard errors in regression and regression-type models.<sup>2</sup>

Let us assume that we wish to make inferences based on a simple random sample drawn from a population. Since we do not observe the whole population, key characteristics like the population mean of y are unknown. We can easily compute the mean and other moments for the sample, and such values are called sample statistics. However, there is no guarantee that the sample statistics are good approximations of the population parameters. Statistical inference is the branch of statistics that is concerned with the problem of gaining knowledge about the values of unknown population parameters using information from sample statistics.

#### 2.1.3.1. Estimation

The term *estimator* refers to the particular method or formula used to obtain sample statistics that are parameter *estimates*. There can be different alternative estimators for a given population parameter. With a few exceptions, different estimators yield distinct estimates of the population parameter of interest.

It is important to note that an estimate itself is a realization of a random variable that follows a probability distribution (or *sampling distribution*). Depending on the particular elements being sampled, sample statistics take on different values. One can view any particular estimate as one of many possible estimates that could have been obtained from multiple, equal-sized random samples drawn from the same population. Thus, the value of the sample mean from a single random sample is only one of numerous sample mean values that could have been calculated from such repeated samples. Moreover, different estimators or estimation methods will often produce different estimates of population parameters, in which case a choice must be made among competing estimators. For example, when the distribution is normal, both the sample median and the sample mean could be used as estimators of the population mean. The sampling distributions of these estimators are different.

Estimators can be judged according to how well they satisfy a few desirable properties. One desirable property of an estimator is unbiasedness. When the expected value of an estimator equals the value of the true parameter being

<sup>2.</sup> Because of this, we do not provide a *notational* distinction between the theoretical response variable and the sampled, or observed, values of the response variable.